

## Journal Pre-proof

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

PII: S2667-1026(23)00021-9  
DOI: <https://doi.org/10.1016/j.imed.2023.03.001>  
Reference: IMED 73



To appear in: *Intelligent Medicine*

Received date: 7 October 2022  
Revised date: 23 February 2023  
Accepted date: 7 March 2023

Please cite this article as: Kun Qian, Ruolan Huang, Zhihao Bao, Yang Tan, Zhonghao Zhao, Mengkai Sun, Bin Hu, Björn W. Schuller, Yoshiharu Yamamoto, Detecting somatisation disorder via speech: introducing the Shenzhen Somatisation Speech Corpus, *Intelligent Medicine* (2023), doi: <https://doi.org/10.1016/j.imed.2023.03.001>

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

© 2023 Published by Elsevier 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/>)

# 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>, and Yoshiharu Yamamoto<sup>7</sup>

1. Key Laboratory of Brain Health Intelligent Evaluation and Intervention,  
Ministry of Education, Beijing Institute of Technology, Beijing 100081, China

2. School of Medical Technology, Beijing Institute of Technology, Beijing 100081, China

3. The First School of Clinical Medicine, Southern Medical University, Guangzhou 510515, China

4. Department of Neurology, Shenzhen University General Hospital, Shenzhen 518055, China

5. GLAM – Group on Language, Audio, & Music, Imperial College London, London SW7 2AZ, UK

6. Chair of Embedded Intelligence for Health Care and Wellbeing, University of Augsburg, Augsburg 86159, Germany

7. Educational Physiology Laboratory, The University of Tokyo, Tokyo 113-0033, Japan

## 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 introduce our somatisation disorder speech database and give benchmark results.

**Methods** By collecting speech samples of somatisation disorder patients, in cooperation with the Shenzhen University General Hospital, we introduce 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 is proposed in our work.

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

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

**Keywords:** Somatisation Disorder, Machine Learning, Healthcare, Computer Audition

## 1. Introduction

According to the World Health Organisation (WHO) web 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 harmful 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 be-

25 behaviours. They may attack people around them or even  
 26 strangers, out of their control. What is more, neurosis  
 27 is also characterised by a higher prevalence in children  
 28 and adolescents [11] and in women than in men [12].  
 29 Compared to depression and anxiety, somatisation dis-  
 30 order (SD) is a mental disorder that is less regarded. In  
 31 fact, it also has a higher prevalence and risk. According  
 32 to a report, the prevalence rate of SD is 10.1 % in Gen-  
 33 eral Hospital Psychiatric Units Tertiary Care Centres in  
 34 India [13]. In addition, the prevalence rate in women is  
 35 higher than that in men. In the study of Babu *et al.*, the  
 36 prevalence rate of SD in adult women was 40.8 % [14].  
 37 SD is a neurosis characterised by persistent fear or be-  
 38 lief in dominance of various somatic symptoms. Be-  
 39 cause SD causes patients to shift from emotional and  
 40 mental distress to the body, they often suffer from un-  
 41 explained physical discomfort [15]. Such physical dis-  
 42 comfort includes stomach pain, back pain, joint pain,  
 43 and further more. This medically unexplained symptom  
 44 is characterised by multiple occurrences, persistence,  
 45 and recurring. Patients are often unaware that they have  
 46 a SD. Attributing this pain to a physical illness, they re-  
 47 peatedly seek clinical medical advice and treatment, but  
 48 to no avail [16]. In terms of economic impact, patients  
 49 frequently seeking medical advice in the wrong direc-  
 50 tion greatly increase their cost of medical care. The per  
 51 capita expenditure for health care of patients is up to  
 52 nine times of the average amount [17]. It also places an  
 53 indirect burden on the healthcare system [18–20]. Ac-  
 54 cording to the report, the annual medical costs of SD in  
 55 the United States are approximately 2 560 billion dol-  
 56 lars [21]. Moreover, unresolved pain and confusion lead  
 57 to lower quality of life and greater mental stress for pa-  
 58 tients, which aggravates their condition, or even cause a  
 59 higher suicide risk [22].

60 However, the vast majority of patients could not re-  
 61 ceive effective treatment. This is because the mental  
 62 health system is severely under-resourced and patients  
 63 lack of relevant mental health knowledge (especially for  
 64 some less common mental disorders). Early diagnosis  
 65 of mental illness is the first step in obtaining beneficial  
 66 treatment for patients, but the approaches are scarce.  
 67 As shown in Fig. 1, clinical scales and speech analy-  
 68 sis are used as the two main methods for the diagnosis  
 69 of mental disorders introduced in this paper. The scale,  
 70 combined with the communication and observation of  
 71 the patient’s behaviour, facial expressions, and speech,  
 72 among others by doctors, is the common form of psy-  
 73 chiatric diagnosis by far. But this approach relies on  
 74 subjective interviews with patients and the clinical expe-  
 75 rience of the doctors, which results in a partial deviation  
 76 of the diagnosis from the reality. With the development

77 of computer technology, it is encouraging to see that,  
 78 artificial intelligence (AI) has been applied to the classi-  
 79 fication of mental disorders, which avoids the pitfalls of  
 80 the scale. In particular, audio signals, because of their  
 81 ‘non-invasive’ nature, combined with the rapidly devel-  
 82 oping computer audition (CA) [23] technology are be-  
 83 coming a popular topic of digital medicine research in  
 84 the search for new digital phenotypes. Speech signals,  
 85 as a subclass of audio signals, have been demonstrated  
 86 to be reliable in the diagnosis of certain mental disor-  
 87 ders, such as depression and anxiety [24–27]. In par-  
 88 ticular, related studies have shown that speech features  
 89 performed better than visual features or text in depres-  
 90 sion prediction tasks [28, 29].

91 There are a mount of AI researches on speech rep-  
 92 resentation of mental disorders. For anxiety, Wang *et al.*  
 93 [24] proposed a new Fourier parameter model us-  
 94 ing the perceptual content of voice quality. Dan *et al.*  
 95 [25] used K-nearest neighbours as classifier and fo-  
 96 cused on the fundamental frequency for classification.  
 97 For depression, much work has been done to detect de-  
 98 pression through speech. Pan *et al.* [26] extracted 988  
 99 speech features from speech data and established a lo-  
 100 gistic regression model to achieve a better depression  
 101 classification rate. Rejaibi *et al.* [27] proposed a deep  
 102 learning-based method to assess depression and pre-  
 103 dict its severity. They aimed to extract Mel Frequency  
 104 Cepstral Coefficients (MFCCs) from speech and used  
 105 long short-term memory networks. For insomnia, Es-  
 106 pinoza *et al.* [30] studied obstructive sleep apnea among  
 107 sleep disorders using a larger speech database contain-  
 108 ing 426 participants. They used supervector or i-vector  
 109 techniques to model speech spectral information and  
 110 predicted by support vector regression. The above re-  
 111 searches demonstrated that speech has the ability to rep-  
 112 resent mental disorders. For SD, unfortunately, due to  
 113 its mix-up feature of psychologically and physically, it  
 114 is challenging for clinicians to recognise this psychiatric  
 115 disorder masquerading as physical pain [31, 32]. Not  
 116 only that, the actual clinical condition may also be di-  
 117 agnosed as a SD, thereby hindering the patient’s real  
 118 need for treatment and longer-term pain [33]. However,  
 119 there are few reports on the application of AI technol-  
 120 ogy in this field. In particular, we found rare work on  
 121 speech as raw data for classification. According to our  
 122 search, Idenfors *et al.* [34] analysed brain images for  
 123 global-brain functional connectivity (GFC). The com-  
 124 bination of GFC values and support vector machine  
 125 (SVM) were used to distinguish patients from the con-  
 126 trols. The results showed that the patients’ GFC was  
 127 abnormal. Lv *et al.* [35] constructed an improved bac-  
 128 terial foraging optimisation-based kernel extreme learn-

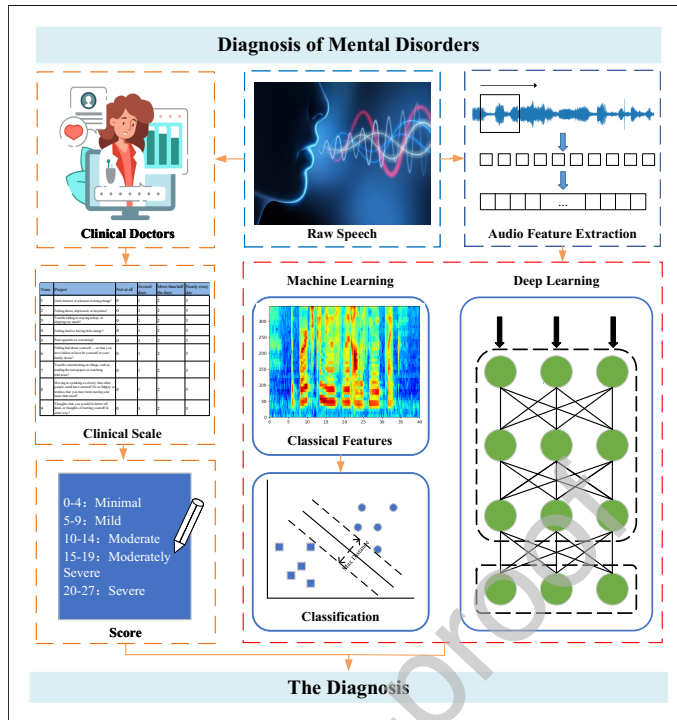


Figure 1: Methods of diagnosing mental disorders.

ing machine (IBFO-KELM) model based on the sym-  
 bol self-assessment scale (SCL-90) data for the diag-  
 nosis of patients with SD. Human communication re-  
 lies 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 informa-  
 tion for diagnosing or predicting mental disorders.

In this paper, we demonstrate the collected data  
 and publish the speech database for the classifica-  
 tion of SD. The name we use is the Shenzhen Som-  
 atisation Speech Corpus (SSSC). Unlike the auto-  
 matic 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 com-  
 posed 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 illu-  
 strate the usability of this database.

The remainder of this paper is organised as follows:  
 At first, the materials and methods are presented in Sec-  
 tion 2. Subsequently, the results of the benchmark work  
 are given in Section 3. We give an experimental dis-  
 cussion, current limitations, and outlooks in Section 4.  
 Finally, we draw a conclusion for our work in Section 5.

## 2. Methods

In this section, we firstly introduce the proposed pub-  
 licly somatisation speech database, i. e., SSSC. Next,  
 OPENSMILE, the toolkit that was used to extract acous-  
 tics 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 par-  
 ticipants involved were informed that their voice data  
 will be used only for research purposes. Their agree-  
 ments for this study were recorded as one of the five

Table 1: Subjects information in data splits.

	Average Age	# Male	# Female
Train	38.2	46	33
Dev	35.0	20	12
Test	36.7	18	12
$\Sigma$	37.1	84	57

176 following original speech phrases. The data was col-  
 177 lected in out-patient-department in the Shenzhen Uni-  
 178 versity General Hospital. We asked the participants to  
 179 speak five sentences (with neutral contextual meaning).  
 180 At the same time, two self-report questionnaires were  
 181 answered by the participants regarding their anxiety and  
 182 SD (physical discomfort). The questionnaire is con-  
 183 structed by two widely used separated questionnaires.  
 184 The two questionnaires are GAD-7 (Patient Anxiety  
 185 Questionnaire) and PHQ-15 (Patient Health Question-  
 186 naire 15) [40–42]. All the data were collected from 12  
 187 November 2020 to 5 April 2021. We used a Shinco  
 188 RV-18 recording pen with 32 GB of storage to record  
 189 all the participants’ voices which have a sample rate of  
 190 32 000 Hz and a bit rate of 16 bps.

191 Following, we give examples of the recorded sen-  
 192 tences which were spoken in Chinese:

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

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

### 204 2.1.2. Data pre-processing

205 As described in [43], we executed a series of data  
 206 pre-processing stages before establishing the ‘standard’  
 207 SSSC, which includes data cleansing, voice activity de-  
 208 tection, speaker diarisation, and speech transcription.  
 209 First, we excluded recordings with low quality (e. g., the  
 210 level of the speech is low compared to the background  
 211 noise). Then, we removed the non-speech parts (e. g.,  
 212 non-subjects’ speech, breathing, and coughing) from  
 213 each recording, which resulted in maintaining only the  
 214 segments including voice from the recordings. The seg-  
 215 ments containing solely the target patient and scripted  
 216 content (e. g., excluding laughing) were kept. Finally,  
 217 we obtained 705 audio recordings from 141 partici-  
 218 pants. To attenuate the effects of the audio record-  
 219 ing equipment, the background noise condition and the  
 220 level of the recording, all files were first high-pass  
 221 filtered to eliminate low-frequency background noise  
 222 (cut-off frequency: 120 Hz, 10<sup>th</sup>-order Chebyshev filter)  
 223 and then their waveforms were normalised individually  
 224 (peak amplitude set to -3 dB).

### 225 2.1.3. Tasks definition

226 We define two tasks for the benchmark setup: First,  
 227 the Anxiety Degree should be grouped into: Yes (la-  
 228 belled as “1” when GAD-7  $\geq$  11), No (labelled as “0”  
 229 when GAD-7 < 11). Then, an estimation of the Physi-  
 230 cal Discomfort Disorder Degree should be made as: Yes  
 231 (labelled as “1” when PHQ-15  $\geq$  10), No (labelled as  
 232 “0” when PHQ-15 < 10).

### 233 2.1.4. Data partitioning

234 Totally, the database contains as mentioned audio  
 235 samples of 705 speech events from 141 subjects who  
 236 were checked without organic disease by experts. The  
 237 number of males in all subjects is 84, and the number  
 238 of females is 57. The average age of all subjects is 37.1  
 239  $\pm$  13.2 years (range 15 to 70). The total duration of  
 240 audio in the database is 3 039.192 s, equalling roughly  
 241 50 minutes. The average sample duration is 4.311  $\pm$   
 242 2.297 s (range 0.864 to 16.920 s). In order to carry out  
 243 the experiment, we randomly partitioned the data into  
 244 a train, a development (dev) and a test set, which are  
 245 subject independence. Table 1 gives more information  
 246 about data partitions in detail. Fig. 3 shows data distri-  
 247 bution details.

## 248 2.2. OPENSMILE and acoustic features

249 OPENSMILE [44] is an open source toolkit which is  
 250 widely applied in the field of acoustic representation ex-  
 251 tractions. OPENSMILE can provide features commonly  
 252 used in classic acoustic signal processing methods, such  
 253 as short-time zero-crossing rate, energy spectrum fea-  
 254 tures, and Mel Frequency Cepstrum Coefficients. To  
 255 get the statistical information of an audio signal sam-  
 256 ple, OPENSMILE firstly extracts the low-level descrip-  
 257 tors (LLDs) from the original frame-level audio signals,  
 258 then performs the statistical information extraction to  
 259 the frame-based LLDs by functionals. By this method,  
 260 the limitations of static machine learning models such  
 261 as SVM are unlocked from the inconsistency of sample  
 262 duration.

263 To train the baseline system, we use the COMPARE  
 264 feature set, which includes 65 LLDs (see Table 2). The

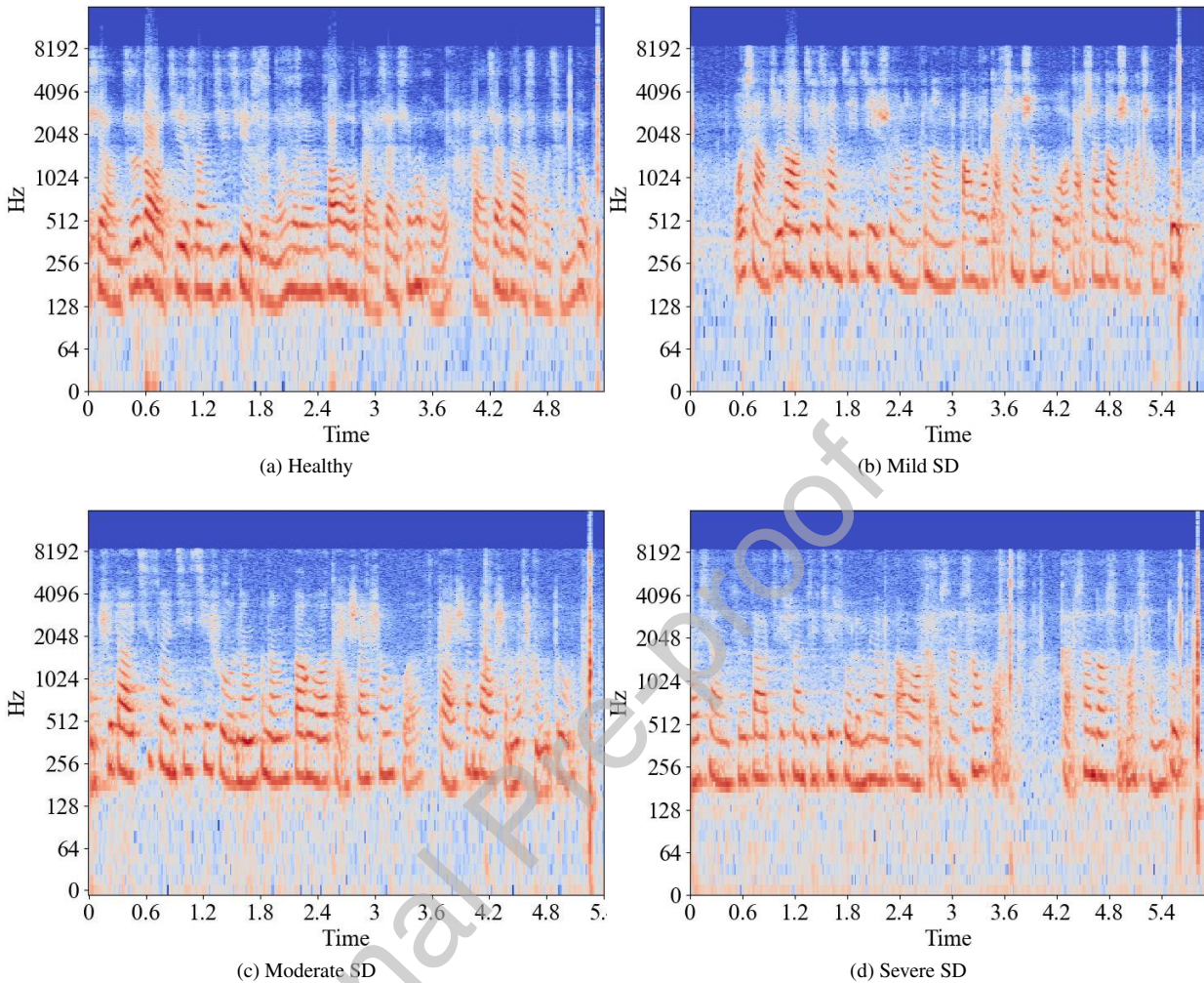


Figure 2: Audio spectrogram examples of different PHQ-15 outcomes.

265 configuration of the 2.3 version OPENSIMILE is COM-  
 266 PARE\_2016. The more specific information and details  
 267 are given in [45]. Moreover, [45] also features a de-  
 268 tailed introduction of the used functionals (see Table 3).  
 269 The COMPARE feature set includes in total 6373 fea-  
 270 tures. The mechanism of functionals is to map the time  
 271 series based LLDs to a scalar value per each applied  
 272 functional (e. g., mean, standard deviation, maximum);  
 273 then, a single, fixed dimension vector which is indepen-  
 274 dent of the audio signal sample's time duration is gen-  
 275 erated [45].

### 276 2.3. Database related questionnaires

277 In our study, we carry out measurement work on two  
 278 scales as mentioned (i. e., GAD-7 and PHQ-15). Based  
 279 on the scores of these two scales, the Shenzhen Uni-

280 versity General Hospital's psychologists enrolled in the  
 281 study classified the conditions into four classes. We give  
 282 a separate introduction of each scale in this part. Specif-  
 283 ically, Table 4 shows the data distribution in details.

284 **GAD-7** Generalised anxiety disorder (GAD) is a  
 285 common and disabling illness that is often underdiag-  
 286 nosed and undertreated [46]. With the influence of the  
 287 COVID-19 in China [47], more and more people are  
 288 suffering from symptoms of anxiety. Therefore, we  
 289 want to get more information of the GAD for further  
 290 research. GAD-7 is a brief clinical measure for as-  
 291 sessing GAD, which consists of 7 items about anxiety  
 292 self-report [48]. The score of each item is between 0-  
 293 3 and the GAD-7's total score is in the range of 0 to  
 294 21. Based on the total scores, the Shenzhen University  
 295 General Hospital's mental health experts classified the

Table 2: The LLDs for COMPARE feature set.

RASTA: Relative Spectral Transform; HNR: Harmonics to Noise Ratio; RMSE: Root Mean Square Energy. Details can be found in [45].

4 Energy related LLDs	Group
RMSE, zero-crossing rate	Prosodic
Sum of auditory spectrum (loudness)	Prosodic
Sum of RASTA-filtered auditory spectrum	Prosodic
6 Voicing related LLDs	Group
Probability of voicing	Voice Quality
$F_0$ (SHS and Viterbi smoothing)	Prosodic
log HNR, jitter (local and $\delta$ ), shimmer (local)	Voice Quality
55 Spectral LLDs	Group
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 (0–8 kHz)	Spectral

296 participants into 4 types (i. e., **no-anxiety** for 0-4, **mild** 324  
 297 **anxiety** for 5-9, **moderate anxiety** for 10-14, and **se-** 325  
 298 **vere anxiety** for 15-21). 326

299 **PHQ-15** SD is a prevalent condition which is not 327  
 300 well treated by many psychiatrists [49]. Patients usu- 328  
 301 ally seek care in the medical setting convinced that they 329  
 302 suffer from physical discomfort rather than a mental dis- 330  
 303 order. Based on this situation, the SD patients often 331  
 304 encounter ineffective care and even harm. That somatising 332  
 305 patients may represent 40 % or more of the ambulatory 333  
 306 medicine patient population greatly magnifies the prob- 334  
 307 lem [50]. To reduce unnecessary expenses, it is very 335  
 308 significant for patients with SD to be fully recognised 336  
 309 and treated. The classification of SD by the speech of 337  
 310 the patients is not only feasible, but also appears effi- 338  
 311 cient and convenient [51]. The PHQ-15 comprises 15  
 312 somatic symptoms from the Patient Health Question-  
 313 naire (PHQ), each symptom scored from 0 (“not both-  
 314 ered at all”) to 2 (“bothered a lot”) [52]. According to  
 315 the PHQ-15 scale’s total score ranging from 0 to 30, the  
 316 participants are divided into 5 types (i. e., **minimal** for  
 317 0-4, **low** for 5-9, **medium** for 10-14, and **high** for 15-  
 318 30).

#### 319 2.4. Machine learning method and optimising strategy

320 **Support Vector Machine** (SVM) is a stable and by 347  
 321 now ‘traditional’ classifier. To make this study compa- 348  
 322 rable and reproducible, we use an SVM classifier with 349  
 323 linear kernel to conduct all experiments. We train an 350

SVM model with the complexity parameter in the range  
 of  $\{10^{-8}, 10^{-7}, \dots, 10^{-1}, 1.0\}$ . Then, we choose the com-  
 plexity 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 acous-  
 tic features, i. e., the full COMPARE feature set, or selec-  
 tively only MFCCs, fundamental frequency ( $F_0$ ), and  
 frequency and bandwidth of the formants ( $F_1$ - $F_3$ ). All  
 the features are extracted by the OPENSMILE feature ex-  
 traction and audio analysis tool.

### 339 3. Results

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

349 In order to obtain a better understanding of the data  
 350 set, we make some changes for the listed labels in Ta-

Table 3: The functionals applied to LLDs in the COMPARE feature set. Note that, the LLDs listed in Table 2 may or may not use all of the functionals of this table, which is described in details in [45].

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

ble 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%. Fig. 4 and Fig. 5 show confusion matrices of the best-performing setup for the test set. From Fig. 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 con-

ditions. Although the classifier performs better on this task, unlike humans, the machine requires more samples to understand the features.

## 4. Discussion

### 4.1. Classification performance

Referring to Table 5 and Table 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.

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.



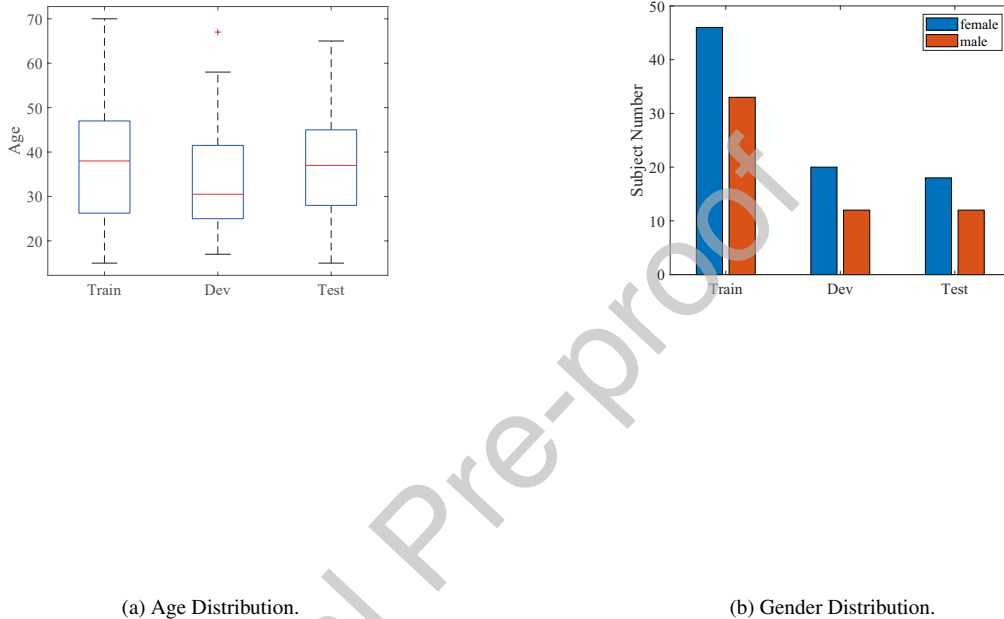


Figure 3: Demographic statistic of the SSSC.

#### 4.2. Features and meaning of SSSC

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 becoming 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.

#### 4.3. Current limitations and outlook

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:

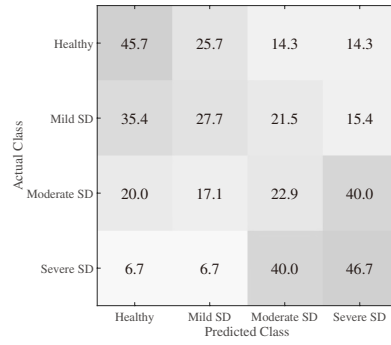
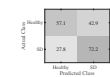
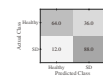


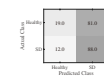
Figure 4: Normalised confusion matrices (values are shown in [%]) of the best-performing (development) model for the PHQ-15 test set.



(a) PHQ-15A



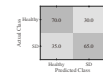
(b) PHQ-15B



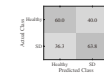
(c) GAD-7&PHQ-15A



(d) GAD-7&PHQ-15B



(e) GAD-7|PHQ-15A



(f) GAD-7|PHQ-15B

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

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

	Score	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

Table 5: Classification results of different feature types and tasks. (Part A)

#: Number of features; 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. (Unit: %.)

Feature type	#	PHQ-15	PHQ-15A	PHQ-15B	GAD-7&PHQ-15A
ComParE	6373	24.4/22.1	53.4/45.9	50.8/49.0	49.8/50.0
MFCCs(all)	1400	28.3/29.5	<b>58.4/64.7</b>	54.8/61.5	51.6/53.5
MFCCs(only coef)	756	26.4/ <b>35.7</b>	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/ <b>54.5</b>

433 1. *Limited dataset.* Due to monocentricity, it 456  
 434 is difficult for us to recruit a sufficient number of 457  
 435 well-represented subjects. This limited amount of data 458  
 436 is not conducive to the application of deep learning 459  
 437 on it. But in fact, that SSSC can be trusted enough 460  
 438 to be adopted for the study of SD was demonstrated 461  
 439 through our benchmark experiments. Therefore, we 462  
 440 encourage researchers to give more consideration to 463  
 441 the application of traditional machine learning which 464  
 442 may be more suitable for the small-scale databases. 465  
 443 More importantly, in future works, data augmentation 466  
 444 is worth looking into and researching, which should be 467  
 445 expected to reduce the model's lack in generalisation. 468  
 446 An easy way is to use some toolboxes for data enhance- 469  
 447 ment. For instance, [53] published a matlab toolkit that 470  
 448 provides 15 different augmentation algorithms for raw 471  
 449 audio data and 8 for spectrograms. 472

451 2. *Overlapping symptoms.* A study has shown over- 474  
 452 lap between somatic symptoms, anxiety, and depres- 475  
 453 sion [54]. In other words, the three are co-morbid and 476  
 454 triggering each other. Thus, different mental illnesses 477  
 455 may show the same features. This may be challenging 478

to model, yet, on the contrary, may facilitate better re-  
 sults. 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 partic-  
 ipants and the clinical experience of doctors. This  
 is an inevitable mistake caused by the lack of objec-  
 tive factors. In the future, one can involve more profes-  
 sional doctors in the scoring, set more differential ver-  
 sions to recognise SD and use multi-modality to ensure  
 reliability of the label. Nevertheless, manual annota-  
 tion 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 intro-  
 ducing 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

Table 6: Classification results of different feature types and tasks. (Part B)

#: Number of features; GAD-7&PHQ15B: 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. (Unit: %.)

Feature type	#	GAD-7&PHQ-15B	GAD-7 PHQ-15A	GAD-7 PHQ-15B
ComParE	6373	51.4/47.6	56.8/50.0	53.9/54.5
MFCCs(all)	1400	52.1/59.2	59.8/53.8	54.3/57.2
MFCCs(only coef)	756	51.2/60.4	<b>61.4</b> /56.2	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</b> /56.8	47.9/50.0	<b>55.3</b> /57.4
F1-F3	14	55.0/ <b>74.8</b>	51.0/ <b>67.5</b>	53.5/ <b>61.9</b>

479 samples to provide researchers with a more scientific  
480 and credible new version.

## 481 5. Conclusion

482 In this study, we firstly introduced a publicly avail-  
483 able speech database, namely SSSC. Then, we de-  
484 scribed the current techniques in somatisation speech  
485 classification. A benchmark experiment was given  
486 based on the approaches proposed in this work. More-  
487 over, we discussed the results and the limitations, and  
488 pointed out some future directions. In this work, an  
489 SVM model trained with features based on the first three  
490 formants, i. e., spectral maxima, performed best in the  
491 classification leading to promising results. We will con-  
492 sider using self-supervision plus fine tuning as strategy  
493 in future work.

## 494 Conflicts of interest statement

495 The authors declare that they have no conflicts of in-  
496 terest.

## 497 Funding

498 This work was partially supported by the Ministry  
499 of Science and Technology of the People’s Republic  
500 of China with the STI2030-Major Projects (Grant  
501 No.2021ZD0201900), the National Natural Sci-  
502 ence Foundation of China (Grant Nos.62227807  
503 and 62272044), the Teli Young Fellow Pro-  
504 gram from the Beijing Institute of Technology,  
505 the Shenzhen Municipal Scheme for Basic Re-  
506 search (Grant Nos.JCYJ20210324100208022 and  
507 JCYJ20190808144005614), China, the JSPS KAK-  
508 ENHI (Grant No.20H00569), the JST Mirai Program  
509 (Grant No.21473074), the JST MOONSHOT Program  
510 (Grant No.JPMJMS229B).

## 511 Author contributions

512 **Kun Qian**: Idea generation, Funding acquisition,  
513 Writing - Conceptualisation, Supervision. **Ruolan**  
514 **Huang**: Data collection & annotation, Funding acqui-  
515 sition, Writing. **Zhihao Bao**: Writing - original draft.  
516 **Yang Tan**: Investigation, Writing. **Zhonghao Zhao**:  
517 Methodology, Writing. **Mengkai Sun**: Writing - re-  
518 view & editing. **Bin Hu**: Writing - review, Supervision,  
519 Funding acquisition. **Björn W. Schuller**: Writing - re-  
520 view. **Yoshiharu Yamamoto**: Writing - review.

## 521 References

- 522 [1] World Health Organisation. Mental disorders. [https://www.who.int/news-room/fact-sheets/detail/](https://www.who.int/news-room/fact-sheets/detail/mental-disorders)  
523 [mental-disorders](https://www.who.int/news-room/fact-sheets/detail/mental-disorders), 2022 (accessed 1 July 2022).  
524  
525 [2] Nina Vindegaard and Michael Eriksen Benros. Covid-19 pan-  
526 demic and mental health consequences: systematic review of  
527 the current evidence. *Brain Behav Immun*, 89:531–542, 2020.  
528 <https://doi.org/10.1016/j.bbi.2020.05.048>.  
529 [3] Nicole Racine, Jessica E Cooke, Rachel Eirich, et al. Child and  
530 adolescent mental illness during covid-19: a rapid review. *Psy-*  
531 *chiatry Res*, 292:113307, 2020. [https://doi.org/10.1016/](https://doi.org/10.1016/j.psychres.2020.113307)  
532 [j.psychres.2020.113307](https://doi.org/10.1016/j.psychres.2020.113307).  
533 [4] Wanderson Carneiro Moreira, Anderson Reis de Sousa, and  
534 Maria do Perpétuo Socorro de Sousa Nóbrega. Mental illness in  
535 the general population and health professionals during covid-19:  
536 a scoping review. *Texto e Contexto Enferm*, 29, 2020. <https://doi.org/10.1590/1980-265x-tce-2020-0215>.  
537 [5] Matthew J Carr, Sarah Steeg, Roger T Webb, et al. Effects of the  
538 covid-19 pandemic on primary care-recorded mental illness and  
539 self-harm episodes in the uk: a population-based cohort study.  
540 *Lancet Public Health*, 6(2):e124–e135, 2021. [https://doi.](https://doi.org/10.1016/S2468-2667(20)30288-7)  
541 [org/10.1016/S2468-2667\(20\)30288-7](https://doi.org/10.1016/S2468-2667(20)30288-7).  
542 [6] World Health Organisation. Mental health and covid-19: early  
543 evidence of the pandemic’s impact: scientific brief, 2 march  
544 2022. [https://www.who.int/publications/i/item/](https://www.who.int/publications/i/item/WHO-2019-nCoV-Sci_Brief-Mental_health-2022.1)  
545 [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),  
546 2022 (accessed 1 July 2022).  
547 [7] World Health Organisation. World health statistics 2022:  
548 monitoring health for the sdgs, sustainable development  
549 goals. [https://www.who.int/publications/i/item/](https://www.who.int/publications/i/item/9789240051157)  
550 [9789240051157](https://www.who.int/publications/i/item/9789240051157), 2022 (accessed 1 July 2022).  
551

- [8] Naomi Launders, Leah Kirsh, David PJ Osborn, 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*, 9(9):725–735, 2022. [https://doi.org/10.1016/S2215-0366\(22\)00225-5](https://doi.org/10.1016/S2215-0366(22)00225-5).
- [9] Joseph Firth, Najma Siddiqi, AI Koyanagi, et al. The lancet psychiatry commission: a blueprint for protecting physical health in people with mental illness. *Lancet Psychiatry*, 6(8):675–712, 2019. [https://doi.org/10.1016/S2215-0366\(19\)30132-4](https://doi.org/10.1016/S2215-0366(19)30132-4).
- [10] Mark Rodgers, Jane Elizabeth Dalton, Melissa Harden, 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*, 18(1), 2018. <https://doi.org/10.5334/ijic.2605>.
- [11] Kathleen Ries Merikangas, Erin F. Nakamura, and Ronald C. Kessler. Epidemiology of mental disorders in children and adolescents. *Dialogues Clin Neurosci*, 11(1):7–20, 2009. <https://doi.org/10.31887/DCNS.2009.11.1/kmerikangas>.
- [12] Walter R Gove. The relationship between sex roles, marital status, and mental illness. *Soc Forces*, 51(1):34–44, 1972. <https://doi.org/10.1093/sf/51.1.34>.
- [13] Kalaivanan Rakesh Chander, Narayana Manjunatha, B Binukumar, et al. The prevalence and its correlates of somatization disorder at a quaternary mental health centre. *Asian J Psychiatr*, 42:24–27, 2019. <https://doi.org/10.1016/j.ajp.2019.03.015>.
- [14] Arjun Rajendra Babu, Alexander John Aswathy Sreedevi, and Vijayakumar Krishnapillai. Prevalence and determinants of somatization and anxiety among adult women in an urban population in kerala. *Indian J Community Med*, 44(Suppl 1):S66–S69, 2019. [https://doi.org/10.4103/ijcm.IJCM\\_55\\_19](https://doi.org/10.4103/ijcm.IJCM_55_19).
- [15] Michael Witthöft, Alexander L Gerlach, and Josef Bäiler. Selective attention, memory bias, and symptom perception in idiopathic environmental intolerance and somatoform disorders. *J Abnorm Psychol*, 115(3):397–407, 2006. <https://doi.org/10.1037/0021-843X.115.3.397>.
- [16] Arthur J Barsky, E John Orav, and David W Bates. Distinctive patterns of medical care utilization in patients who somatize. *Med Care*, 44(9):803–811, 2006. <https://doi.org/10.1097/01.mlr.0000228028.07069.59>.
- [17] G Richard Smith Jr, Roberta A Monson, and Debby C Ray. Psychiatric consultation in somatization disorder. *N Engl J Med*, 314(22):1407–1413, 1986. <https://doi.org/10.1056/NEJM198605293142203>.
- [18] Frauke Dorothee Weiss, Winfried Rief, and Maria Kleinstäuber. Health care utilization in outpatients with somatoform disorders: descriptives, interdiagnostic differences, and potential mediating factors. *Gen Hosp Psychiatry*, 44:22–29, 2017. <https://doi.org/10.1016/j.genhosppsych.2016.10.0>.
- [19] Zbigniew Jerzy Lipowski et al. Somatization: the concept and its clinical application. *Am J Psychiatry*, 145(11):1358–1368, 1988. <https://doi.org/10.1176/ajp.145.11.1358>.
- [20] Per Fink. The use of hospitalizations by persistent somatizing patients. *Psychol Med*, 22(1):173–180, 1992. <https://doi.org/10.1017/S0033291700032827>.
- [21] Arthur J Barsky, E John Orav, and David W Bates. Somatization increases medical utilization and costs independent of psychiatric and medical comorbidity. *Arch Gen Psychiatry*, 62(8):903–910, 2005. <https://doi.org/10.1001/archpsyc.62.8.903>.
- [22] Andrea P Chioqueta and Tore C Stiles. Suicide risk in patients with somatization disorder. *Crisis*, 25(1):3, 2004. <https://doi.org/10.1027/0227-5910.25.1.3>.
- [23] Kun Qian, Xiao Li, Haifeng Li, et al. Computer audition for healthcare: opportunities and challenges. *Front Digit Health*, 2:5, 2020. <https://doi.org/10.3389/fdgh.2020.00005>.
- [24] Kunxia Wang, Ning An, Bing Nan Li, et al. Speech emotion recognition using fourier parameters. *IEEE Trans Affect Comput*, 6(1):69–75, 2015. <https://doi.org/10.1109/TAFFC.2015.2392101>.
- [25] Marius Dan Zbancioc and Silvia Monica Feraru. A study about the automatic recognition of the anxiety emotional state using emo-db. In *2015 E-Health and Bioengineering Conference (EHB)*, pages 1–4, Iasi, Romania, 2015. IEEE. <https://doi.org/10.1109/EHB.2015.7391506>.
- [26] Wei Pan, Jingying Wang, Tianli Liu, et al. Depression recognition based on speech analysis. *Chin Sci Bull*, 63(20):2081–2092, 2018. <https://doi.org/10.1360/N972017-01250>.
- [27] Emna Rejaibi, Ali Komaty, Fabrice Meriaudeau, et al. Mfcc-based recurrent neural network for automatic clinical depression recognition and assessment from speech. *Biomed Signal Process Control*, 71:103–107, 2022. <https://doi.org/10.1016/j.bspc.2021.103107>.
- [28] James R Williamson, Elizabeth Godoy, Miriam Cha, et al. Detecting depression using vocal, facial and semantic communication cues. In *Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge*, pages 11–18, Amsterdam, The Netherlands, 2016. <https://doi.org/10.1145/2988257.2988263>.
- [29] Le Yang, Hichem Sahli, Xiaohan Xia, et al. Hybrid depression classification and estimation from audio video and text information. In *Proceedings of the 7th annual workshop on audio/visual emotion challenge*, pages 45–51, New York, the United States, 2017. <https://doi.org/10.1145/3133944.3133950>.
- [30] Fernando Espinoza-Cuadros, Rubén Fernández-Pozo, Doroteo T Toledano, et al. Reviewing the connection between speech and obstructive sleep apnea. *Biomed Eng Online*, 15(1):1–20, 2016. <https://doi.org/10.1186/s12938-016-0138-5>.
- [31] Robert C Smith and Francesca C Dwamena. Classification and diagnosis of patients with medically unexplained symptoms. *J Gen Intern Med*, 22(5):685–691, 2007. <https://doi.org/10.1007/s11606-006-0067-2>.
- [32] Alessia Raffagnato, Caterina Angelico, Perla Valentini, et al. Using the body when there are no words for feelings: alexithymia and somatization in self-harming adolescents. *Front Psychiatry*, 11:262–261, 2020. <https://doi.org/10.3389/fpsy.2020.00262>.
- [33] L De Jonge, S Petrykiv, J Fennema, et al. Misdiagnosis of loin pain hematuria syndrome as a somatization disorder. *Eur Psychiatry*, 41(S1):S491–S491, 2017. <https://doi.org/10.1016/j.eurpsy.2017.01.597>.
- [34] Pan Pan, Yangpan Ou, Qinji Su, et al. Voxel-based global-brain functional connectivity alterations in first-episode drug-naive patients with somatization disorder. *J Affect Disord*, 254:82–89, 2019. <https://doi.org/10.1016/j.jad.2019.04.099>.
- [35] Xinen Lv, Huiling Chen, Qian Zhang, et al. An improved bacterial-foraging optimization-based machine learning framework for predicting the severity of somatization disorder. *Algorithms*, 11(2):17–34, 2018. <https://doi.org/10.3390/a11020017>.
- [36] Cynthia B Solot, Debbie Sell, Anne Mayne, et al. Speech-language disorders in 22q11. 2 deletion syndrome: best practices for diagnosis and management. *Am J Speech Lang Pathol*, 28(3):984–999, 2019. [https://doi.org/10.1044/2019\\_AJSLP-16-0147](https://doi.org/10.1044/2019_AJSLP-16-0147).
- [37] Arthur Brito-Marcelino, Edmea Fontes Oliva-Costa, Salvyana Carla Palmeira Sarmento, et al. Burnout syndrome in speech-

- language pathologists and audiologists: a review. *Rev Bras Med Trab*, 18(2):217, 2020. <https://doi.org/10.47626/1679-4435-2020-480>.
- [38] Alyssa M Lanzi, James M Ellison, and Matthew L Cohen. 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*, 6(5):987–1002, 2021. [https://doi.org/10.1044/2021\\_PERSP-20-00295](https://doi.org/10.1044/2021_PERSP-20-00295).
- [39] Heejin Kim, Mark Hasegawa-Johnson, Adrienne Perlman, Jon Gunderson, Thomas S Huang, Kenneth Watkin, and Simone Frame. Dysarthric speech database for universal access research. In *Ninth Annual Conference of the International Speech Communication Association*, 2008.
- [40] Yen-Cheng Shih, Chien-Chen Chou, Yi-Jiun Lu, et al. Reliability and validity of the traditional chinese version of the gad-7 in taiwanese patients with epilepsy. *J Formos Med Assoc*, pages 1–7, 2022. <https://doi.org/10.1016/j.jfma.2022.04.018>.
- [41] Francisco Cano-García, Javier, Roger Muñoz-Navarro, Albert Abad, Sesé, et al. Latent structure and factor invariance of somatic symptoms in the patient health questionnaire (phq-15). *J Affect Disord*, 261:21–29, 2020. <https://doi.org/10.1016/j.jad.2019.09.077>.
- [42] Xiao-Jie Huang, Hai-Yan Ma, Xue-Mei Wang, 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*, 311:327–335, 2022. <https://doi.org/10.1016/j.jad.2022.05.079>.
- [43] Jing Han, Kun Qian, Meishu Song, et al. An early study on intelligent analysis of speech under covid-19: severity, sleep quality, fatigue, and anxiety. In *Proceedings of INTERSPEECH*, pages 4946–4950, Shanghai, China, 2020. <https://doi.org/10.1109/EHB.2015.7391506>.
- [44] Florian Eyben, Martin Wöllmer, and Björn Schuller. Opensmile: the munich versatile and fast open-source audio feature extractor. In *Proceedings of the 18th ACM International Conference on Multimedia*, pages 1459–1462, Firenze, Italy, 2010. <https://doi.org/10.1145/1873951.1874246>.
- [45] Florian Eyben. *Real-time Speech and Music Classification by Large Audio Feature Space Extraction*. Springer International Publishing, Cham, Switzerland, 2015. Doctoral Thesis. <https://doi.org/10.1007/978-3-319-27299-3>.
- [46] Jeremy DeMartini, Gayatri Patel, and Tonya L Fancher. Generalized anxiety disorder. *Ann Intern Med*, 170(7):ITC49–ITC64, 2019. <https://doi.org/10.7326/AITC201904020>.
- [47] Yeen Huang and Ning Zhao. Generalized anxiety disorder, depressive symptoms and sleep quality during covid-19 outbreak in china: a web-based cross-sectional survey. *Psychiatry Res*, 288:112954, 2020. <https://doi.org/10.1016/j.psychres.2020.112954>.
- [48] Robert L Spitzer, Kurt Kroenke, Janet BW Williams, et al. A brief measure for assessing generalized anxiety disorder: the gad-7. *Arch Intern Med*, 166(10):1092–1097, 2006. <https://doi.org/10.1001/archinte.166.10.1092>.
- [49] François Mai. Somatization disorder: a practical review. *Can J Psychiatry*, 49(10):652–662, 2004. <https://doi.org/10.1177/070674370404901002>.
- [50] Robert C Smith. Somatization disorder. *J Gen Intern Med*, 6(2):168–175, 1991. <https://doi.org/10.1007/BF02598318>.
- [51] Thomas E Oxman, Stanley D Rosenberg, Paula P Schnurr, et al. Linguistic dimensions of affect and thought in somatization disorder. *Am J Psychiatry*, 1985. <https://doi.org/10.1176/ajp.142.10.1150>.
- [52] Kurt Kroenke, Robert L Spitzer, and Janet BW Williams. The phq-15: validity of a new measure for evaluating the severity of somatic symptoms. *Psychosom Med*, 64(2):258–266, 2002. <https://doi.org/10.1097/0006842-200203000-00008>.
- [53] Gianluca Maguolo, Michelangelo Paci, Loris Nanni, and Ludovico Bonan. Audiogmenter: a matlab toolbox for audio data augmentation. *Applied Computing and Informatics*, 2021. <https://doi.org/10.1108/ACI-03-2021-0064>.
- [54] Xiaoya Fu, Fengyu Zhang, Feng Liu, et al. Brain and somatization symptoms in psychiatric disorders. *Front Psychiatry*, 10:146, 2019. <https://doi.org/10.3389/fpsy.2019.00146>.
- [55] Weiwei Liu, Haobo Wang, Xiaobo Shen, et al. The emerging trends of multi-label learning. *IEEE Trans Pattern Anal Mach Intell*, 44(11):7955–7974, 2021. <https://doi.org/10.1109/TPAMI.2021.3119334>.
- [56] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. *Advances in neural information processing systems*, 33:22243–22255, 2020.
- [57] Hankook Lee, Sung Ju Hwang, and Jinwoo Shin. Self-supervised label augmentation via input transformations. In *International Conference on Machine Learning*, pages 5714–5724. PMLR, 2020. <https://doi.org/10.1109/TPAMI.2021.3119334>.
- [58] Kun Qian, Zixing Zhang, Alice Baird, and Björn Schuller. Active learning for bird sound classification via a kernel-based extreme learning machine. *The Journal of the Acoustical Society of America*, 142(4):1796–1804, 2017. <https://doi.org/10.1121/1.5004570>.