

# Interactive robot with multimodal multitask model for early screening of multiple common adolescent mental disorders

Peiwu Qin

[pwqin@sz.tsinghua.edu.cn](mailto:pwqin@sz.tsinghua.edu.cn)

Tsinghua University <https://orcid.org/0000-0002-7829-8973>

Zhicheng Du

Likun Zhang

Shiyao Zhai

Zhengyang Lei

Yongjie Zhou

Yu Dongmei

Chenggang Yan

Hangzhou Dianzi University

Xi Yuan

Hangzhou Dianzi University

Jiansong Ji

Yang Liu

Zhenglin Chen

---

## Article

**Keywords:** Adolescent mental disorder, Mental health screening, Interactive multi-sensor robot, Multimodal learning, Human-Computer interaction, Computer-aided screening

**Posted Date:** March 25th, 2025

**DOI:** <https://doi.org/10.21203/rs.3.rs-5731226/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

**Additional Declarations:** There is **NO** Competing Interest.

---

1 **Title**

2 Interactive robot with multimodal multitask model for early screening of multiple  
3 common adolescent mental disorders

4  
5 **Authors**

6 Zhicheng Du<sup>1,2</sup>, Yang Liu<sup>2</sup>, Likun Zhang<sup>1,2</sup>, Shiyao Zhai<sup>1,2</sup>, Zhengyang Lei<sup>1,2</sup>, Yongjie  
7 Zhou<sup>3</sup>, Dongmei Yu<sup>4</sup>, Chenggang Yan<sup>5</sup>, Xi Yuan<sup>5</sup>, Zhenglin Chen<sup>1,2,6,\*</sup>, Jiansong Ji<sup>6,\*</sup>,  
8 Peiwu Qin<sup>1,2,6\*</sup>

9 **Affiliations**

10 <sup>1</sup>Center of Precision Medicine and Healthcare, Tsinghua-Berkeley Shenzhen Institute,  
11 Shenzhen, Guangdong Province, 518055, China

12 <sup>2</sup>Institute of Biopharmaceutics and Health Engineering, Tsinghua Shenzhen  
13 International Graduate School, Shenzhen, Guangdong Province, 518055, China

14 <sup>3</sup>Shenzhen Mental Health Center, Shenzhen Kangning Hospital, Shenzhen, Guangdong  
15 Province, 518055, China

16 <sup>4</sup>School of Mechanical, Electrical & Information Engineering, Shandong University,  
17 Weihai, Shandong Province, 264209, China

18 <sup>5</sup>School of Communication Engineering, Hangzhou Dianzi University, Hangzhou,  
19 Zhejiang Province, 310018, China

20 <sup>6</sup>Zhejiang Key Laboratory of Imaging and Interventional Medicine, Department of  
21 Radiology, Lishui Central Hospital, The Fifth Aliated Hospital of Wenzhou Medical  
22 University

23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  

---

\* Corresponding authors: [pwqin@sz.tsinghua.edu.cn](mailto:pwqin@sz.tsinghua.edu.cn), [jjstcty@wmu.edu.cn](mailto:jjstcty@wmu.edu.cn),  
[chenzlin1992@163.com](mailto:chenzlin1992@163.com).

41 **Abstract**

42 The early detection of mental disorders in adolescents represents a significant global  
43 public health challenge. Due to the complex and subtle nature of mental disorders,  
44 making it difficult to detect abnormalities using a single factor. Additionally, the  
45 generalized multimodal Computer-Aided Screening (CAS) systems, incorporating  
46 interactive robots for adolescent mental health assessment, remain unavailable. In this  
47 study, we present an Android application equipped with mini-games and chat recording,  
48 deployed in a portable robot, to screen 3,783 middle school students. This system  
49 generates a multimodal screening dataset comprising facial images, physiological  
50 signals, voice recordings, and textual transcripts. We develop a model called **GAME**  
51 (**G**eneralized Model with Attention and **M**ultimodal **E**mbraceNet) with novel attention  
52 mechanism that integrates cross-modal features into the model. GAME evaluates  
53 adolescent mental conditions with high accuracy (73.34% – 92.77%) and F1-Score  
54 (71.32% – 91.06%) and outperforms traditional methods. Our findings reveal that each  
55 modality contributes dynamically to mental disorder detection and the identification of  
56 comorbidities across various disorders, supporting the feasibility of an explainable  
57 model. This study provides a system capable of acquiring multimodal information and  
58 constructs a generalized multimodal integration algorithm with novel attention  
59 mechanisms for the early screening of adolescent mental disorders.

60

61 **Keywords:** Adolescent mental disorder, Mental health screening, Interactive multi-  
62 sensor robot, Multimodal learning, Human-Computer interaction, Computer-aided  
63 screening.

64

65

66 **Main**

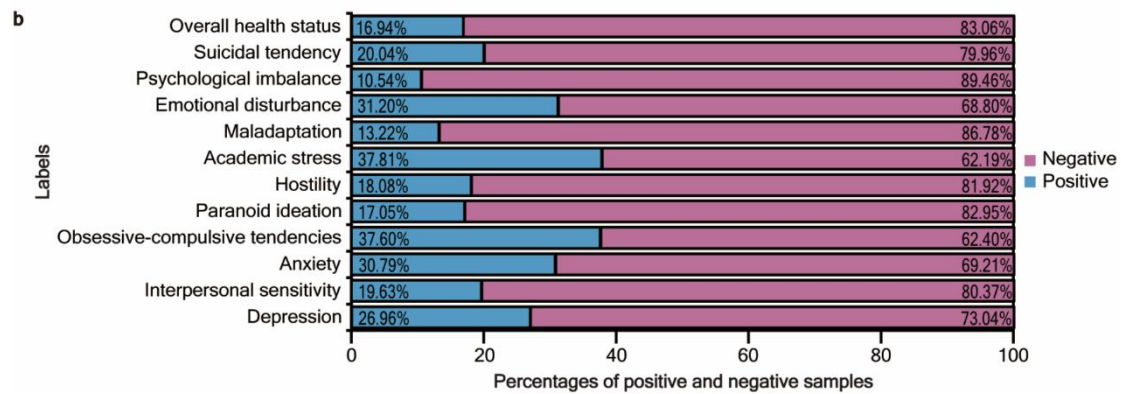
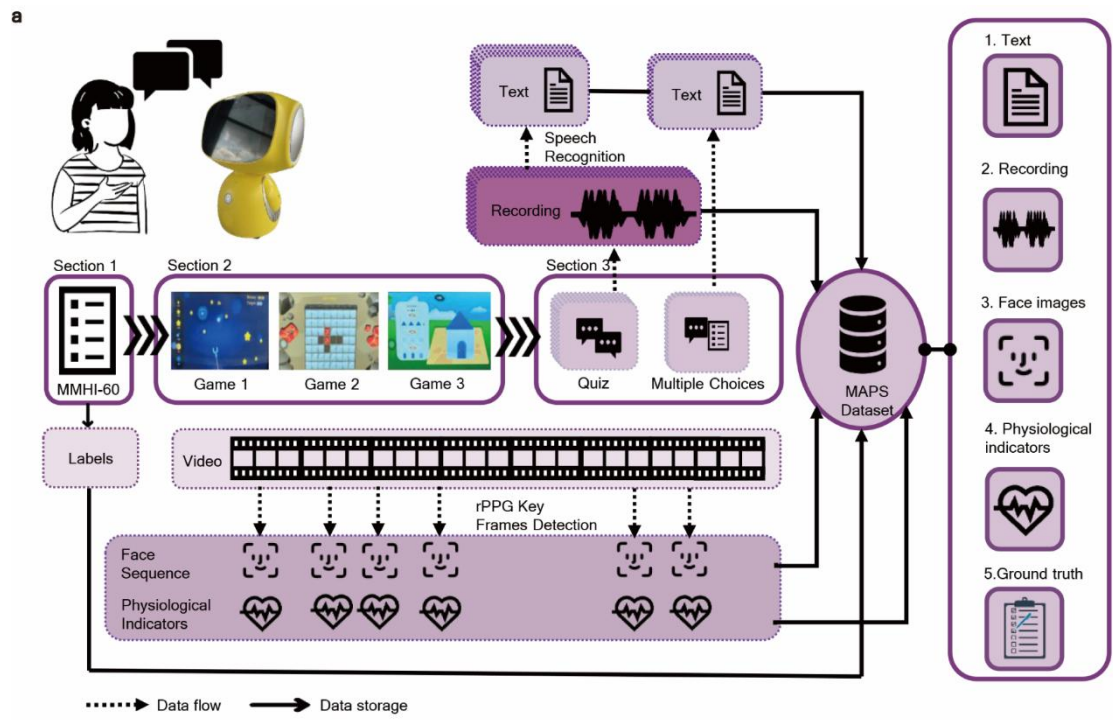
67 Adolescence is a crucial period of life development during which significant  
68 psychosocial adjustments takes place. A large percentage of mental health disorders that  
69 progress into adulthood exhibit symptoms at a young age<sup>1,2</sup>, indicating that adolescent  
70 mental health issues could degenerate into worse later-life illnesses. Approximately 13%  
71 of adolescents aged 10–19 in the world are diagnosed with different types of mental  
72 illness, of which 80 million adolescents aged 10–14 and 86 million adolescents aged  
73 15–19 are deeply affected by mental disorders<sup>3,4</sup>. Unfortunately, ~80% adolescents are  
74 unable to receive precise and professional psychological counseling when they demand  
75 mental health services<sup>5</sup> and ~50% adolescents with mental disorders have access to  
76 psychotherapy<sup>6</sup>. Traditional screening methods for mental disorders include  
77 questionnaires and interviews<sup>7</sup>, where the results rely on patients' self-reports and  
78 psychiatrists' observations<sup>8,9</sup>. However, these methods are inherently susceptible to  
79 subjective bias. Furthermore, barriers like stigma in disclosing mental illness or  
80 negative attitudes towards professionals<sup>10</sup> lead to inaccurate psychological assessments  
81 and a vicious cycle of disease deterioration. To address these limitations, interactive  
82 robots providing an enjoyable and acceptable interface with less defensive altitude and  
83 hostility offer a promising avenue for unconscious screening<sup>11</sup>. The humanoid robot is  
84 more accurate at detecting pediatric mental health problems than parental or child self-

85 reporting<sup>12</sup>. Therefore, imperceptible and interactive screening robot with  
86 corresponding algorithm for accurate and opportune screening to adolescent mental  
87 disorders can support healthcare agencies and ameliorate the social burden<sup>13,14</sup>.

88 Here, we develop a humanoid robot equipped with well-designed emotional stimuli  
89 that facilitates the acquisition of the Multimodal Adolescent Psychological Screening  
90 (MAPS) dataset (age 12–15), including facial images, physiological indicators, audio  
91 recordings, and textual transcripts (**Fig. 1**). The acquired multimodal dataset is analyzed  
92 with statistical model to minimize the distance between prediction and ground-truth  
93 provided by screening questionnaires. The Mental Health Inventory of Middle School  
94 Students (MMHI-60)<sup>15,16</sup> is a screening questionnaire specially designed to assess  
95 Chinese adolescents' mental health and has exhibited high specificity and sensitivity in  
96 screening 10 different types of mental disorders (Supplementary **Methods**). We  
97 maintain MMHI-60 questionnaire to screen 10 types of mental disorders with additional  
98 screening results suggested by experienced psychologists for suicidal tendency. Thus,  
99 a total of 12 psychological conditions are labeled as ground truth for individual subject  
100 in the dataset, including: (1) depression, (2) interpersonal sensitivity, (3) anxiety, (4)  
101 obsessive-compulsive tendencies, (5) paranoid ideation, (6) hostility, (7) academic  
102 stress, (8) maladaptation, (9) emotional disturbance, (10) psychological imbalance, (11)  
103 suicidal tendency, and (12) overall mental health status<sup>17</sup>.

104 Robotic platforms with human-computer interaction have been utilized for  
105 intervention in adolescent mental health<sup>18-20</sup>. However, existing systems lack a  
106 computer-aided screening (CAS) algorithm for psychometrics, The CAS approach has  
107 shown promise in diagnosing of mental disorders in adolescents<sup>21,22</sup>, which can process  
108 different types of input data (e.g., physical activity, sociability, device usage patterns,  
109 etc.) collected from various sensors<sup>23</sup> are utilized to recognize specific mental disorders  
110 including depression, anxiety, and stress<sup>24-29</sup>. However, current CAS models employing  
111 single-modal feature encounter limitations in constructing a comprehensive  
112 representation of the latent multimodal feature space<sup>30</sup>, which weakens their  
113 performance. Multimodal CAS models have been used to predict psychological  
114 disorders and mental states by feature importance ranking, feature selection, and feature  
115 concatenation strategies<sup>31-34</sup>. Nevertheless, the screening of specific psychiatric  
116 disorders and the lack of interpretability of these models have hindered the adoption of  
117 CAS models in clinical applications. Limited exploration exists on whether a  
118 generalized model with interpretability could accurately screen adolescents' mental  
119 disorders. Therefore, achieving both generalization and interpretability in the CAS  
120 system remains a challenge for clinical utility.

121



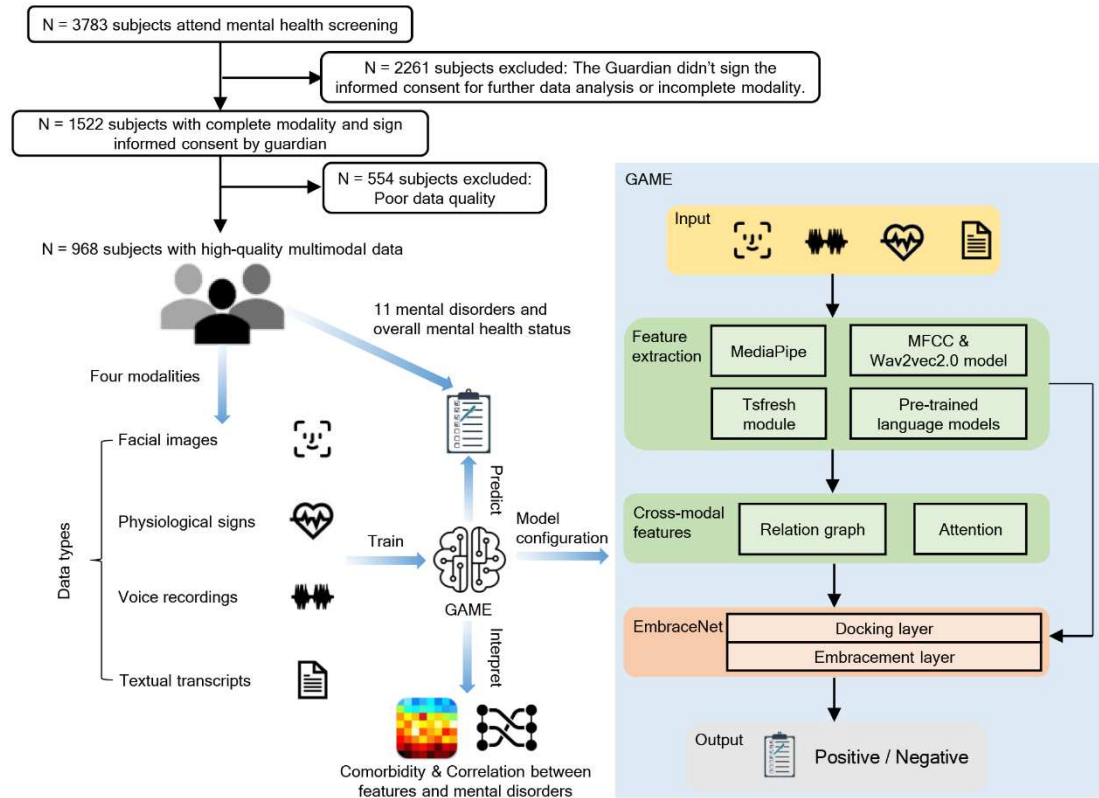
122  
 123 **Figure 1. MAPS data acquisition and database construction.** **a**, The flowchart of data acquisition.  
 124 A humanoid robot with a customer-designed Android application that can interact with participants  
 125 is used for data collection. The data collection procedure has three consecutive sections:  
 126 psychological screening, emotional stimuli games, and questions-and-answers. The mental health  
 127 inventory is designed by psychological expert as labels of ground truth. Using remote photo-  
 128 plethysmography (rPPG) and available processing algorithm, key frames (i.e., images with clear  
 129 and unmasked face) and physiological indicators are extracted from videos captured during the  
 130 games and questions-and-answers sections. The responses in the questions-and-answers sections  
 131 are recorded and converted into text using the speech recognition technique of Iflyrec  
 132 (<https://www.iflyrec.com>), supplied by iFlytek. **b**, the sample distributions for 11 types of screened  
 133 mental disorders and overall mental health status. The ratios depict the imbalance of the MAPS  
 134 dataset, and the positive samples labeled as ‘Overall mental health status’ represent the abnormal  
 135 adolescents.

136  
 137 Hence, we propose GAME (Generalized model with Attention and Multimodal  
 138 EmbraceNet), a generalized model based on distance-weighted attention mechanisms  
 139 and multimodal feature fusion in the EmbraceNet backbone network<sup>35</sup> (**Fig. 2**) for

140 adolescent mental disorders screening. GAME extracts eight single-modal features  
141 named Expression, Expression nuance, and Eye movement from face images;  
142 Physiological signs; MFCC and Wav2vec from audio recordings; PERT and RoBERTa  
143 from textual transcripts, respectively. Inspired by the diagnostic strategies employed by  
144 psychologists during structured diagnostic and screening interviews with adolescents<sup>36</sup>,  
145 we propose a novel attention mechanism for multi-scale feature to integrate inter-model  
146 correlation weights and eight single-modal features. Additionally, we introduce cross-  
147 modal features named Relation graph and Attention, which play a crucial role in extract  
148 deeper information and alleviate the interference of noisy features. Hyper-Emotion  
149 theory<sup>37,38</sup> indicates that adolescents suffered from mental disorders have abnormal  
150 multimodal emotional and behavioral responses to the same interactive stimuli in  
151 contrast to healthy subjects. GAME, guided by the Hyper-Emotion theory, accurately  
152 predicts overall mental health status and identifies 11 types of adolescent mental  
153 disorders based on multimodal responses. We harness GAME's capabilities to predict  
154 comorbidities among adolescents with multiple mental disorders and compare the  
155 findings with relevant studies. The ablation experiment that involve the stepwise  
156 removal of individual modal inputs and fusion analyses to evaluate contribution ratio  
157 of each model from trained GAME. These experiments collectively affirm the  
158 significance of modal features and the robustness of multimodal fusion within our  
159 framework.

160 In summary, this study develops a cost-effective and highly precise screening robot  
161 platform along with GAME to screen early mental illness among adolescents. The  
162 development of a practical and adolescent-friendly mental health screening system,  
163 tailored to adolescents and capable of delivering accurate and interpretable results,  
164 holds significant promise for the integration of CAS systems within clinical contexts.  
165 The theory-consistent comorbidity prediction underscores the GAME's reliability for  
166 predicting comorbidity from data-driven perspective. GAME excels in identifies the  
167 dominated features for certain mental disorder and provides valuable guidance in the  
168 design of screening protocols, especially when dealing with single-modal data. This  
169 guidance recommends the clinician prioritizes critical features and directs researchers  
170 towards uncovering implicit patterns or theories through a data perspective.

171



172

173 **Figure 2. Pipeline of data processing and GAME's structure.** A total of 3,787 people participated  
 174 in the mental health screening, retaining 968 samples after exclusion. Based on four types of input,  
 175 GAME has been trained to predict mental disorders, mining comorbidity and correlation between  
 176 multimodal features and mental disorders in adolescent. MediaPipe, Mel-  
 177 Frequency Cepstrum Coefficients (MFCC), Wav2vec2.0, Tsfresh module, pre-trained language  
 178 models including Robustly Optimized BERT approach (RoBERTa), Pre-  
 179 training BERT with Permuted Language Model (PERT) are used to extract single-modal features  
 180 from facial images, voice recording, physiological indicators, and textual transcripts respectively.  
 181 The extracted features undergo task-level fusion, and then two cross-modal features are generated  
 182 through unimodal features. Eight single-modal and two cross-modal features are fused by  
 183 EmbraceNet. BERT means Bidirectional Encoder Representations from Transformers.

184

## 185 Results

### 186 Multimodal database construction

187 We construct MAPS dataset with 3,787 Chinese middle school students aged 12 to 15  
 188 and filter to 968 (**Fig. 2** and Supplementary **Method**). This dataset spans across four  
 189 distinct data modalities, encapsulating a spectrum of 11 mental disorders and overall  
 190 mental health status. The 12 mental health conditions in the dataset have different  
 191 distribution and the imbalanced positive-to-negative ratios (**Fig. 1b**), which are ranked  
 192 from high to low as follows: obsessive-compulsive tendencies (6.56), interpersonal  
 193 sensitivity (5.31), overall mental health status (4.90), academic stress (4.87), hostility  
 194 (4.53), psychological imbalance (4.09), suicidal tendency (2.71), depression (2.44),  
 195 emotional disturbance (2.25), anxiety (2.21), maladaptation (1.66), paranoid ideation  
 196 (1.64). The subjects are hailing from diverse multi-centers and cities within Guangdong

197 Province China. The MAPS dataset collects comprehensive features via portable  
 198 screening platform compared to the public mental disorder dataset. The IMAGEN  
 199 study<sup>39</sup> and the Adolescent Brain Cognitive Development Study (ABCD)<sup>40</sup> are large  
 200 multimodal adolescent mental health datasets, which encompass diverse modalities  
 201 such as MRI neuroimaging and behavioral assessments. There are also private clinical  
 202 datasets that have been used to train AI models for the diagnosis of specific adolescent  
 203 psychiatric disorders. However, the current datasets are not compatible with portable  
 204 screening for mental disorders due to data privacy, high cost constraints and intricate  
 205 data acquisition processes. MAPS uses a readily accessible and inexpensive data  
 206 collection platform, facilitating seamless scalability for large-scale population  
 207 screening (**Supplementary Table 2**).

208

### 209 **Attention mechanism and multimodal integration**

210 With extracted single-modal and cross-modal features, we compare reported machine  
 211 learning (ML) models used for mental disorders diagnosis<sup>41-43</sup>, including Support  
 212 Vector Machine with Polynomial Kernel (SVM-Poly) and Radial Basis Function  
 213 (SVM-RBF) Kernel, Random Forest (RF), and Gradient-Boosting Decision Tree  
 214 (GBDT) with GAME, to evaluate the prediction accuracy for 12 mental conditions and  
 215 robustness of GAME. The assessment criteria for these models are predicated on  
 216 accuracy and weighted F1-Score, bolstered by 10-fold stratified cross-validation  
 217 methodology instead of the random split to evaluate the model's performance. GAME  
 218 averagely enhances the accuracy of 3.31% - 76.24% (SVM-RBF), 3.31% - 76.55%  
 219 (SVM-Poly), 3.31% - 15.49% (RF), and 3.93% - 17.98% (GBDT) in comparison to the  
 220 bracket's baseline models (**Table 1**). In terms of model robustness, GAME enhances  
 221 the weighted F1-score of the SVM-RBF, SVM-Poly, RF, and GBDT models by 5.07%  
 222 - 83.31%, 6.57% - 83.94%, 6.34% - 23.78%, and 6.08% - 22.87%, respectively. The  
 223 wide-ranging improvements indicates the efficacy of GAME in mental disorders  
 224 screening.

225

226 **Table 1 | Models evaluation and comparison for 12 different prediction tasks**

Ground truth	Evaluation metric	SVM-RBF	SVM-Poly	RF	GBDT	GAME
Overall mental health status	Accuracy	70.15% (84.30%, 49.17%)	64.99% (83.06%, 30.68%)	83.08% (84.40%, 81.61%)	82.23% (83.68%, 80.99%)	89.26% (92.78%, 87.63%)
	F1-Score	69.78% (79.47%, 52.87%)	64.70% (79.49%, 31.32%)	76.82% (79.21%, 75.37%)	76.86% (79.09%, 74.84%)	87.49% (91.92%, 85.42%)
Depression	Accuracy	60.98% (74.38%, 27.38%)	59.38% (74.38%, 31.50%)	72.86% (73.45%, 71.59%)	71.30% (73.14%, 69.94%)	80.16% (82.47%, 78.13%)
	F1-Score	56.87% (66.04%, 12.49%)	55.29% (66.97%, 16.68%)	63.35% (65.12%, 61.89%)	64.19% (66.32%, 61.61%)	76.80% (79.15%, 74.00%)



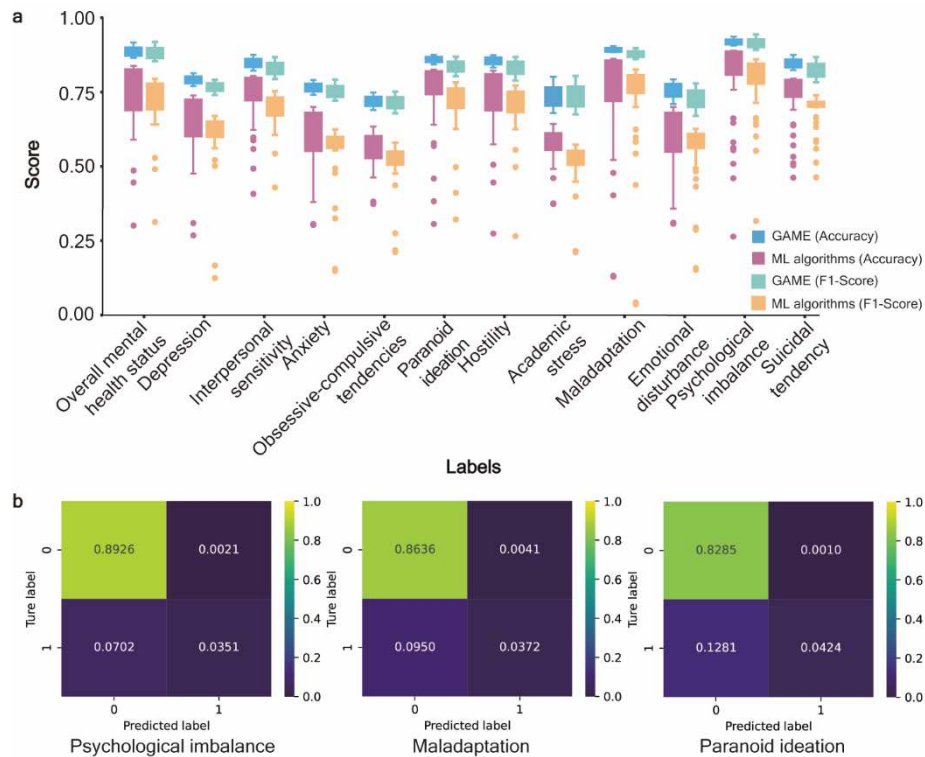
Interpersonal sensitivity	Accuracy	70.29% (80.99%, 56.42%)	66.16% (80.37%, 41.31%)	80.15% (80.58%, 79.03%)	79.07% (80.27%, 78.41%)	85.85% (88.66%, 83.33%)
	F1-Score	68.56% (75.28%, 60.59%)	64.68% (74.86%, 42.94%)	72.59% (74.26%, 71.63%)	72.88% (74.43%, 71.12%)	82.76% (86.72%, 79.37%)
Anxiety	Accuracy	57.04% (70.46%, 31.20%)	54.08% (70.56%, 30.89%)	68.38% (68.91%, 66.84%)	66.44% (67.98%, 64.78%)	77.58% (80.21%, 75.26%)
	F1-Score	52.01% (61.63%, 15.53%)	49.38% (62.42%, 14.89%)	58.21% (61.40%, 56.47%)	59.49% (61.81%, 56.71%)	74.83% (79.18%, 72.08%)
Obsessive-compulsive tendencies	Accuracy	53.67% (63.95%, 38.22%)	51.68% (62.60%, 37.91%)	60.87% (62.71%, 57.65%)	59.25% (61.68%, 55.90%)	73.04% (76.04%, 70.10%)
	F1-Score	49.44% (58.00%, 21.87%)	46.41% (55.40%, 21.21%)	52.50% (56.97%, 48.60%)	54.89% (57.14%, 51.63%)	71.32% (75.17%, 67.85%)
Paranoid ideation	Accuracy	72.17% (82.95%, 46.48%)	64.69% (82.96%, 31.20%)	82.58% (83.06%, 81.91%)	81.38% (82.33%, 80.57%)	87.08% (88.66%, 85.42%)
	F1-Score	70.13% (78.28%, 49.79%)	63.74% (75.82%, 32.21%)	75.71% (76.59%, 75.18%)	75.74% (76.85%, 74.48)	83.59% (86.92%, 80.33%)
Hostility	Accuracy	70.95% (82.74%, 51.13%)	64.81% (82.13%, 28.00%)	81.47% (81.92%, 80.37%)	80.41% (81.20%, 79.34%)	86.78% (88.54%, 84.38%)
	F1-Score	69.21% (77.25%, 55.57%)	63.66% (76.93%, 26.59%)	74.24% (75.82%, 73.37%)	74.50% (75.90%, 72.78%)	83.54% (86.78%, 78.88%)
Academic stress	Accuracy	57.10% (64.88%, 38.12%)	54.67% (64.57%, 37.91%)	61.11% (63.64%, 58.69%)	59.53% (61.99%, 56.20%)	74.18% (81.25%, 69.07%)
	F1-Score	49.52% (57.32%, 21.58%)	47.82% (56.60%, 21.15%)	53.02% (56.13%, 49.27%)	54.98% (56.90%, 50.19%)	73.06% (80.46%, 67.48%)
Maladaptation	Accuracy	68.16% (86.36%, 13.84%)	62.79% (86.78%, 13.53%)	86.30% (86.78%, 85.22%)	84.83% (85.33%, 83.99%)	90.08% (91.67%, 89.58%)
	F1-Score	67.17% (82.58%, 4.33%)	62.77% (80.64%, 3.71%)	80.67% (81.31%, 79.98%)	80.32% (80.92%, 79.83%)	87.65% (89.77%, 86.16%)
Emotional disturbance	Accuracy	55.67% (70.35%,	53.06% (68.91%,	68.74% (70.56%,	66.61% (69.11%,	77.17% (80.41%,

		31.61%)	31.30%)	67.56%)	63.85%)	72.16%)
	F1-Score	50.56% (62.67%, 15.88%)	48.23% (62.51%, 15.23%)	59.22% (62.43%, 56.09%)	59.70% (62.36%, 56.34%)	73.00% (77.92%, 67.01%)
Psychological imbalance	Accuracy	75.15% (89.46%, 51.44%)	70.49% (89.46%, 26.96%)	89.25% (89.46%, 88.22%)	88.25% (88.84%, 86.88%)	92.77% (94.79%, 91.75%)
	F1-Score	76.19% (85.99%, 60.09%)	71.65% (84.49%, 31.66%)	84.44% (84.57%, 84.13%)	84.43% (84.98%, 83.30%)	91.06% (94.38%, 89.18%)
Suicidal tendency	Accuracy	68.45% (80.06%, 46.77%)	69.46% (79.96%, 50.91%)	79.53% (79.96%, 78.20%)	78.20% (79.34%, 76.24%)	85.43% (88.66%, 83.51%)
	F1-Score	65.71% (73.90%, 46.34%)	66.44% (72.58%, 51.30%)	71.27% (71.81%, 70.85%)	71.70% (73.48%, 70.20%)	82.20% (86.72%, 78.27%)

227 The outcomes of ML algorithms are the average values of single-modal features and cross-modal  
228 features, while the outputs of GAME are the average values assessed by the 10-fold stratified  
229 cross-validation method. Data in red denotes the highest value in the row, while data in blue  
230 denotes the row's next-highest value. The maximum and minimum values are denoted by the two-  
231 tuple results in parentheses.

232

233 Specially, we integrate the baseline outcomes of ML algorithms (**Fig. 3**) to  
234 juxtapose them with the GAME concerning their predictive efficacy across diverse  
235 manifestations of mental disorders. The results shows that GAME enhances accuracy  
236 by 5.8% - 52.78% (Depression), 4.86% - 44.54% (Interpersonal sensitivity), 7.02% -  
237 46.70% (Anxiety), 9.09% - 35.13% (Obsession-compulsive tendencies), 4.03% - 55.89%  
238 (Paranoid ideation), 4.03% - 58.78% (Hostility), 9.30% - 36.27% (Academic stress),  
239 3.93% - 76.55% (Maladaptation), 6.61% - 45.87% (Emotional disturbance), 3.31% -  
240 65.81% (Psychological imbalance), 5.37% - 38.66% (Suicidal tendency), and 4.86% -  
241 58.58% (Overall mental health status), while the weighted F1-Score of GAME is  
242 boosted by 10.92% - 64.31% (Depression), 7.49% - 39.82% (Interpersonal sen sitivity),  
243 13.68% - 59.94% (Anxiety), 18.53% - 50.11% (Obsessive-compulsive tendencies),  
244 7.95% - 51.39% (Paranoid ideation), 6.29% - 56.95% (Hostility), 19.12% - 51.91%  
245 (Academic stress), 5.07% - 83.94% (Maladaptation), 11.75% - 57.77% (Emotional  
246 disturbance), 6.57% - 59.40% (Psychological imbalance), 10.54% - 35.86% (Suicidal  
247 tendency), and 8.28% - 56.17% (Overall mental health status), respectively.  
248 Furthermore, we employ the metrics of weighted precision, weighted recall, and the  
249 normalized confusion matrix to rigorously evaluate the performance of GAME across  
250 several classification tasks (Supplementary **Fig. 10-12**). GAME outperforms ML  
251 methods in both binary and multiple classification indicated by various metrics.



252

253 **Figure 3. Evaluation results of comparison between GAME and ML algorithms in various**  
 254 **mental disorders. a**, the results assessed by the accuracy and weighted F1-score in order to evaluate  
 255 the performance of GAME and ML algorithms work in predicting various types of mental disorders,  
 256 while the values of ML algorithms are incorporated in accordance with those distinct types of mental  
 257 disorders. **b**, the top three mental conditions predicted by GAME, which are assessed by normalized  
 258 confusion matrix in 10-fold stratified cross-validation.

259

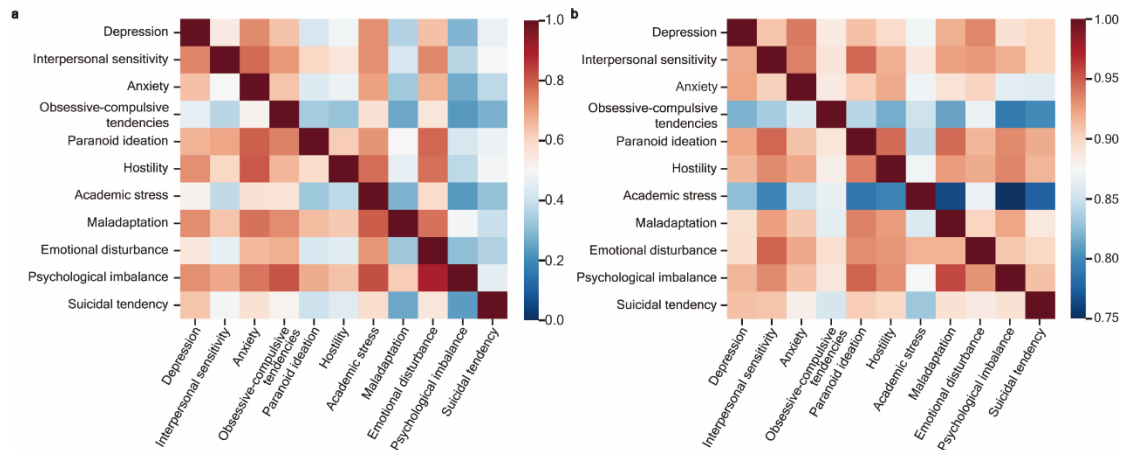
260 **Comorbidity among various mental disorders**

261 We use correlation analysis to evaluate the comorbidities and relevancy levels  
 262 among different mental disorders in adolescents (**Fig. 4a**). The findings indicate that:  
 263 (1) A significant comorbidity exists between depression and anxiety in young  
 264 individuals; (2) Adolescents with anxiety exhibit an elevated susceptibility to emotional  
 265 disturbances; (3) Adolescents who suffer from depression and anxiety tend to  
 266 experience heightened levels of academic stress; (4) Adolescents with interpersonal  
 267 sensitivity disorder manifest an increased vulnerability to emotional disturbance,  
 268 anxiety, depression, and academic stress, where anxiety and depression are more  
 269 prevalent; (5) Teenagers with paranoid ideation are more susceptible to anxiety,  
 270 obsessive-compulsive tendencies, and emotional disturbance; (6) Hostility and  
 271 maladaptation are associated with higher levels of academic stress and psychological  
 272 imbalance. Also, a correlation between hostility and anxiety is discernible. (7) Within  
 273 our cohort displaying psychological imbalances, we note a high occurrence of  
 274 emotional disturbance, followed by academic stress and obsessive-compulsive  
 275 tendencies. (8) Suicidal tendencies in adolescents may be influenced more easily by  
 276 depression, anxiety, academic stress, and emotional disturbance. (Detail analysis is  
 277 shown in Supplementary **Results**). Co-morbidities or correlations among different

278 mental disorders aligns with the findings presented in existing published literature and  
 279 clinical reports, thus reinforcing the validity of our data-driven approaches in reaching  
 280 concordant conclusions with clinical evidence.

281 In addition, we observe novel comorbidities via the prediction ability of GAME  
 282 (Fig. 4b). The potential comorbidities are inferred from GAME prediction but are not  
 283 revealed by correlation analysis, for example: (1) maladaptation and paranoid ideation  
 284 are closely linked to psychological imbalance; (2) there is a comorbidity between  
 285 paranoid ideation and hostility as well as maladaptation; (3) there is a comorbidity  
 286 between suicidal tendency with interpersonal sensitivity and paranoid ideation; (4)  
 287 emotional disturbance has a comorbidity with interpersonal sensitivity. (Further details  
 288 in the Supplementary Results). A quantitative measure of the comorbidity between  
 289 different mental disorders and complex interactions can be estimated with our method.  
 290 The attention mechanism in this study employs the dual relationship in calculating the  
 291 feature distance, which can be extended to multiple feature similarities when more data  
 292 points are available later.

293



294

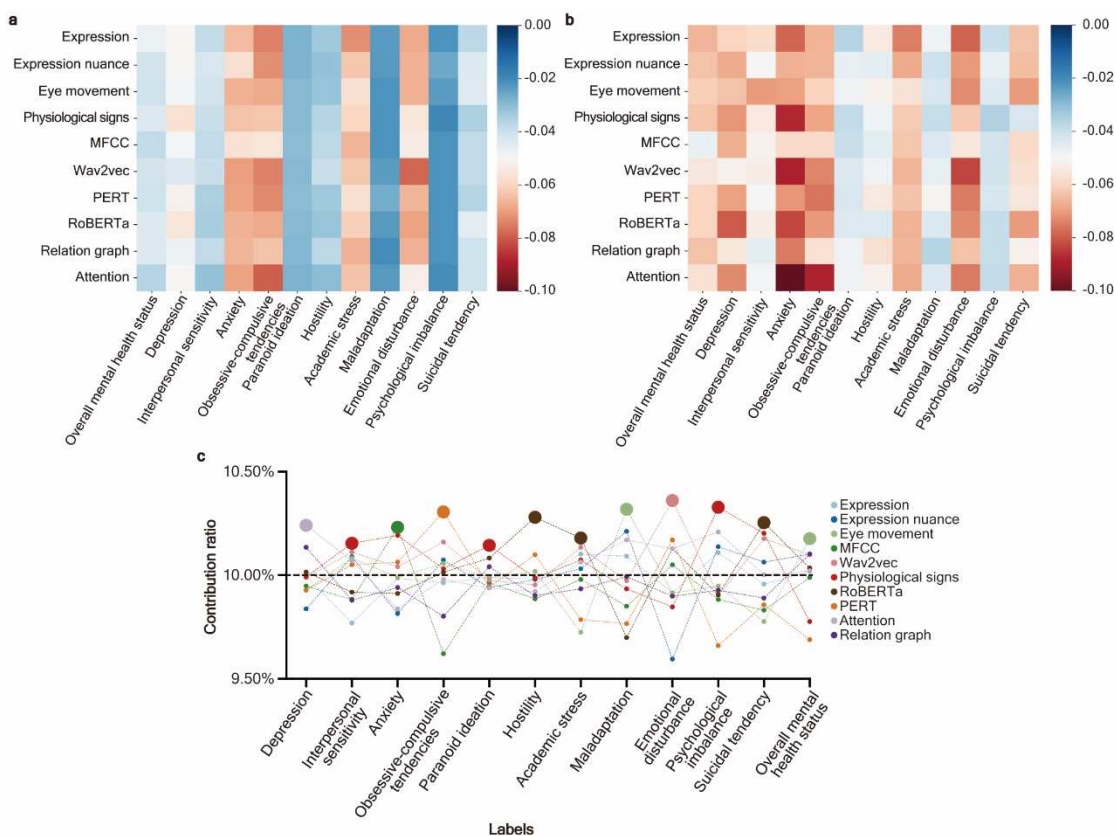
295 **Figure 4. Comorbidities among 11 different mental disorders in adolescents.** **a**, the heat map  
 296 reports the comorbidity association through data statistics. The value of color bar indicates the  
 297 correlation ratio, which are calculated by the number of samples who are simultaneously suffering  
 298 from two different mental disorders. **b**, the heat map shows the correlation of GAME predictions.  
 299 The score of color bars is calculated based on the accuracy obtained from GAME with various  
 300 model parameters trained by different mental disorders data, with higher accuracy indicating greater  
 301 resemblance between the two mental disorder. Darker blue indicates poorer correlation while deeper  
 302 red indicates higher correlation.

303

### 304 **Modality ablation experiments**

305 Each modal feature can boost the GAME's accuracy in predicting various mental  
 306 disorder (Fig. 5a). The impact of different modal features on the performance of GAME  
 307 varies, with some exerting more pronounced effects than others, which facilitates  
 308 GAME's ability to explain the specific contributions of each modality to the prediction  
 309 of particular mental disorders. The modal features can be ranked based on their  
 310 contribution to the model's accuracy, with the following order from highest to lowest:  
 311 Wav2vec, Expression, RoBERTa, Expression nuance, Relation graph, Eye movement,

312 PERT, Attention, Physiological signs, and MFCC. The absence of specific modal  
 313 features can result in a considerable decline in the prediction accuracy of GAME when  
 314 predicting specific mental disorders, such as Attention features and obsessive-  
 315 compulsive tendencies, Wav2vec features and emotional disturbance, expression  
 316 features and academic stress. In terms of weighted F1-score (**Fig. 5b**), the average  
 317 contribution of modal features to the robustness and stability of GAME is listed in  
 318 descending order: Attention, RoBERTa, Expression, PERT, Eye movement, Wav2vec,  
 319 Expression nuance, Physiological signs, Relation graph, and MFCC. Analogously, the  
 320 removal of certain modal features can greatly diminish the robustness of GAME; for  
 321 example, Expressions, Physiological signs, Wav2vec, Roberta, and Attention facilitate  
 322 GAME’s stability in predicting anxiety. In addition, Attention and Wav2vec help  
 323 GAME improve accuracy and robustness in the tasks of screening obsessive-  
 324 compulsive tendencies and emotional disturbance. The results explainably demonstrate  
 325 the hierarchical importance of various factors in mental disorder prediction.  
 326



327  
 328 **Figure 5. Ablation and contribution ratio for different modal features.** **a**, the heat map shows  
 329 the impact of modal feature elimination on prediction accuracy of GAME. The score of color bar  
 330 indicates the percentage of accuracy decrease and the symbol ‘-’ represents decline. **b**, the influence  
 331 on weighted F1-Score after removing certain modal feature of GAME. Deeper red denotes better  
 332 correlation, while darker blue suggests lower correlation. **c**, the line chart describes the contribution  
 333 ratio of different features in various GAME prediction tasks, which provides the interpretation of  
 334 the reasoning why GAME provides this screening decision.

335  
 336 **Modal feature contributions**

337 GAME indicates the dynamic contribution of each modal feature throughout the  
338 multimodal feature fusion to tailor the needs of different scenarios, underscoring the  
339 adaptability of modal features in predicting different mental disorders (**Fig. 5c**). This  
340 analysis establishes associations between specific mental disorders and their most  
341 significant diagnostic features, including Attention and Depression; Physiological signs  
342 and Interpersonal sensitivity; MFCC (i.e., voice recording) and Anxiety; PERT (i.e.,  
343 textual transcripts) and Obsession-compulsive tendencies; Physiological signs and  
344 Paranoid ideation; RoBERTa (i.e., textual transcripts) and Hostility; RoBERTa and  
345 Academic stress; Eye movement and Maladaptation; Wav2vec (i.e., voice recording)  
346 and Emotional disturbance; Physiological signs and Psychological imbalance;  
347 RoBERTa and Suicidal tendency; as well as Eye movement and Overall mental health  
348 status. These findings explain the deterministic features utilized by GAME to make  
349 predictions for certain mental disorders, which are consistent with the screening  
350 methods used in previous work<sup>44,45</sup> (Detailed analysis in Supplementary **Results**). In  
351 resource- or time-limiting scenarios, the conclusion about important feature provides  
352 guidance for choosing the most valuable modality for certain mental disorder screening,  
353 thus optimizing the efficiency of mental disorder.

354

## 355 **Discussion**

356 CAS models for biomedical applications have experienced rapid development<sup>46-48</sup> and  
357 multimodal learning increasingly gaining traction in the domain of disease screening  
358 and diagnosis<sup>49,50</sup>. Nevertheless, the absence of screening hardware slows down the  
359 proliferation of CAS within the psychological sphere, subsequently limits the creation  
360 of a generalized and interpretable multimodal CAS for screening adolescent mental  
361 disorders. Addressing this, we design and create an interactive robot with a well-  
362 designed Android APP to screen adolescent disorders unobtrusively across a broad  
363 population. Then we build the MAPS database and develop a generalized multimodal  
364 model, named as GAME, which showcases commendable accuracy and robustness in  
365 predicting adolescent mental ailments. The integration of multiple feedback features is  
366 a promising predictor of psychological disorders in adolescents.

367 The multimodal feature fusion and the incorporation of attention mechanism boost  
368 the universality of GAME in the task of screening diverse mental disorders, where  
369 previous deep learning models are developed specifically for certain mental  
370 disorders<sup>51,52</sup>. GAME evaluates adolescent's mental health conditions with an accuracy  
371 of 73.34% – 92.77%, a F1-Score of 71.32% – 91.06%, a specificity of 73.24% – 93.14%  
372 and a sensitivity of 73.04% – 92.77%. Since other psychometric tools were reported to  
373 have ~70% specificity<sup>53,54</sup>, GAME is a more effective and powerful tool for screening  
374 adolescent mental disorders. Modality ablation shows that each modal feature provides  
375 a positive contribution in predicting performance. Notably, the absence of Attention  
376 leads to a ~10% reduction in model performance when predicting anxiety and  
377 obsession-compulsive tendencies. In a nutshell, GAME is superior to conventional ML  
378 algorithms and screening tools in prediction performance due to its thorough feature  
379 extraction and cross-modal information mining.

380 Comorbidity is not a rarity<sup>55</sup>, emphasizing the importance of comprehensive

381 analyses for a detailed psychological profiling of adolescents. Adolescents with mental  
382 disorders require comorbidity analysis to create a precise psychological portrait.  
383 Comorbidities hold profound clinical implications for the diagnosis of mental disorders,  
384 the prescription of appropriate treatments, and the long-term management<sup>56</sup>. However,  
385 to the best of our knowledge, few researchers utilize multimodal algorithms to mine  
386 comorbidities among adolescent psychological disorders. GAME can quantify the  
387 relevancy magnitude between different mental disorders in adolescents, which  
388 improves the accuracy of the mental disorders screening and provides insights for  
389 development of adolescent psychological theories through data-driven perspective. For  
390 example, GAME predicts a comorbidity between emotional disturbance and  
391 interpersonal sensitivity, shown in empirical research<sup>57</sup>, which indicates that unstable  
392 social relationships cause emotional disorders. GAME as a digital assistant to prompt  
393 the psychiatrist to give priority to the interpersonal sensitivity rather than emotional  
394 disturbance. The GAME can be extended to discover novel comorbidities if more modal  
395 features and mental disorder types are provided.

396 Interpretability is crucial for the development and application of CAS systems in  
397 clinical settings. Unexplained or opaque models (known as "black boxes") make it  
398 difficult to understand the logic reasoning of clinical decision<sup>58</sup>. By dissecting the  
399 trained GAME's parameters, we explain how GAME makes predictions through the  
400 contribution ratio for each modal feature during diverse prediction tasks, which  
401 demystifies the intricate interplay between mental disorders and modal features through  
402 modality ablation. For example, GAME suggests that Physiological signs is more  
403 important than other modal features in predicting interpersonal sensitivity, which is  
404 consistent with the report that interpersonal sensitivity is associated with higher systolic  
405 blood pressure<sup>59</sup>. GAME guides future research directions through comorbidity  
406 relationships and correlation between features and mental disorders. For instance,  
407 GAME predicts that maladjustment and paranoid ideation are possibly linked to  
408 psychological imbalance. However, there is currently no relevant work to show the  
409 comorbidity between them, and future work is required to fill this gap.

410 This study is not without its limitations. First, even that GAME has been validated,  
411 the size of the MAPS dataset is modest, which restricts the performance of data-driven  
412 models and necessitates the collection of larger samples to enable GAME to learn subtle  
413 features about adolescent mental disorders. Adolescents' mental disorders are closely  
414 related to their living environment<sup>60</sup>. In the future, we can enlarge the MAPS dataset to  
415 include more cities and countries with diverse economical stages, geographical  
416 environments, and social culture. Second, the materials of emotional stimuli may not  
417 be abundant enough. To improve the reliability of audiovisual stimuli<sup>61,62</sup>, emotionally  
418 elicited film clips should be included. Third, public multimodal datasets can be used to  
419 train GAME for widespread applications. However, multimodal datasets for screening  
420 of adolescent mental disorder are not available. Transfer learning with a pre-trained  
421 model can be adopted to extra psychometric applications instead of screening. Fourth,  
422 GAME can be extended to tackle the issue of modalities absence, which has not been  
423 addressed in computational psychology. Real-world datasets often contain inadequate  
424 modality data for a variety of reasons, like data privacy, failed acquisitions, data

425 corruption, and costly testing<sup>63</sup>. The missing modality problem has been studied in other  
426 diseases' diagnosis<sup>64</sup>.

427 In summary, this study elucidates that an economically viable (< \$400), portable,  
428 interactive, expandable robot with vivid emotional stimulation materials can effectively  
429 facilitate screening and diagnosis of adolescent mental health disorders. GAME,  
430 underpinned by robust theoretical frameworks, has the advantages of high accuracy,  
431 strong stability, and interpretability, which presents a promising avenue in the realm of  
432 mental disorder screening and unveil the relationship among various mental disorders  
433 as well as the correlation between mental disorders and modalities from a model-driven  
434 perspective.

435

## 436 **Methods**

437 Approval for the study was granted by the Office of Research Ethics at Tsinghua  
438 University, Shenzhen International Graduate School under Protocol No. 41 in 2021.

439

## 440 **Design of Android application**

441 The Android application's architecture encompasses data transfer and database  
442 management, built upon a foundation of technological components including: Spring  
443 Boost 2.0, Spring Cloud, MySQL, VUE, Docker, Remote Dictionary Server (REDIS),  
444 and QUEUE technologies, etc. The development process consists of two distinct  
445 phases: protocol design and code implementation. Firstly, we collaborate and consult  
446 with professional psychologists, psychological counselors from middle school, and  
447 representative parents to identify the requirements and appropriate tools for adolescent  
448 mental disorders screening. Subsequently, we formulate the interaction scheme and  
449 functional architecture of the application. Once we validate the engineering feasibility  
450 of the scheme and structure, we proceed with designing the user interface (UI) and user  
451 experience (UE). We follow the code development order of application (APP) client,  
452 application programming interface (API) server, and background database management  
453 system. In detail, we use Java and the front-end framework VUE for development of  
454 the application client, employ Restful API and Domain-driven Design (DDD)  
455 technologies for application API server development, and utilize REDIS and MySQL  
456 for background database management systems. Upon completing the application  
457 development, we conduct application program testing, including App content testing,  
458 App performance testing, App function testing, App visual testing, debugging, and  
459 repairing bugs. Finally, we deploy the application onto the interactive robot for on-site  
460 screening (Supplementary Fig. 1–9). The screening platform we develop provides  
461 objective and involuntary screening appropriate for repetitive screening, addressing the  
462 bias associated with questionnaire-based screening. Moreover, the APP's content  
463 facilitates personalized further development, allowing researchers to tailor different  
464 stimulus materials to meet the various demands of psychological screening and  
465 diagnosis.

466

## 467 **MAPS Dataset Collection**

468 Our adolescent multimodal mental health screening dataset contains facial, textual,



469 acoustic, and physiological data, four data modalities, which are collected from  
470 multiple middle schools in Guangdong Province with 3,783 volunteers ranging from 12  
471 to 15 years old and filtered to 968 after exclusion (Supplementary **Methods**). Each data  
472 is collected by a humanoid robot. The main components of this robot include a touch  
473 screen, a camera, a speaker, and a recording device. The touch screen displays the test  
474 content and allows interaction with the test taker. The camera records video of the  
475 volunteers' faces, and the recording device records the volunteers' voices during the test.  
476 The recorded data is transferred to a configured personal computer for storage. An  
477 Android app installed in the robot system completes the entire testing and data  
478 collection process (Supplementary **Methods**). Personal information, such as gender,  
479 age, class number, and student ID, is required prior to data collection. The volunteer  
480 will enter all of the above information into the robot via the touch screen. The recorded  
481 video of the acquisition process and classroom environment is provided in the  
482 Supplementary **Videos** and Supplementary **Fig. 13**.

483 To minimize the physical and psychological discomfort experienced by adolescent  
484 participants during screening caused by a wearable device, we use a high-resolution  
485 camera installed into the robot to collect video data and calculate physiological signs  
486 by the rPPG algorithm integrated in the back-end server. The rPPG<sup>65</sup> algorithm, coined  
487 as non-contact PPG<sup>66,67</sup>, is a technique to analyze the face video to extract physiological  
488 indicators, including heart rate, heart rate variability, changes in blood pressure, and  
489 respiration rate. Stress and relaxation levels can be calculated using a DL algorithm and  
490 the arousal-valence emotion model<sup>68,69</sup> based on physiological indicators. Eventually,  
491 we obtain six physiological metrics and save them in the database. The volunteer may  
492 move significantly during the screening process, potentially causing the rPPG  
493 algorithm to fail at deriving certain physiological indicators. Only the key and clear  
494 frames in the videos identified by the rPPG algorithm can be used to acquire the  
495 physiological indicators, and we save the pairs of face images and physiological signs  
496 to maintain a consistent correspondence between them.

#### 497 **MMHI-60**

498 The MMHI-60 is adapted from the Symptom Checklist-90 (SCL-90)<sup>70</sup>, which was  
499 designed through a two-year follow-up survey on the mental problems of middle school  
500 students in more than 100 schools across China and has been successfully applied to  
501 the mental disorders screening for Chinese middle school students<sup>71</sup>. The MMHI-60  
502 comprises 60 questions to measure relevant symptoms of 10 distinct mental problems  
503 (including depression, interpersonal sensitivity, anxiety, obsessive-compulsive  
504 tendencies, paranoid ideation, hostility, academic stress, maladaptation, emotional  
505 disturbance, and psychological imbalance). For each question, the respondent assigns a  
506 score ranging from 1 to 5, depending on whether they have recently undergone a  
507 specific type of symptom or behavior, which represents none, mild, moderate, heavy,  
508 and serious, respectively<sup>72</sup>. The MMHI-60 uses a 5-point Likert scale, where a score of  
509 2-2.99 indicates the presence of mild problematic symptoms; 3-3.99 suggests moderate  
510 symptoms; 4-4.99 indicates the presence of severe symptoms; and a rating of 5 denotes  
511 severe psychological symptoms. Final score is the average score of its corresponding  
512

513 questions, allowing the participants to be identified as having the potential for  
514 symptoms of a relative mental disorder. The mental health issue is recognized when the  
515 average score of the subscale is equal to or higher than 2, which will be regarded as  
516 positive. The ground truth of overall mental health status is obtained by combining all  
517 the scores from subscales (i.e., the higher the score, the worse the overall mental health  
518 status), and the ground truth of suicidal tendency is obtained by both the MMHI-60 and  
519 diagnostic advice from the psychiatrist. The question list of the MMHI-60 is presented  
520 in the Supplementary **Methods**.

### 521 **Theoretical Supporting framework**

522 This work relies on Hyper-Emotion theory, which supports GAME a theoretical  
523 foundation for the plausibility of predicting psychological conditions based on the  
524 magnitude of emotional responses to external stimuli within adolescents. It posits that  
525 mental diseases stem from a cognitive appraisal process that undergoes a series of  
526 unconscious transitions culminating in the manifestation of fundamental emotions, such  
527 as happiness or anger. The Hyper-Emotion theory contains five principles: (1) The  
528 principle of unconscious transitions to fundamental emotions. People develop a series  
529 of unconsciously shifts from a physiological sensation or cognitive assessment to a  
530 fundamental emotion that are contextually appropriate to the circumstance but aberrant  
531 in its response intensity. Such transitions lead to the start of a psychological illness, but  
532 they persist during the illness<sup>38</sup>. (2) The principle of no voluntary control. Individuals  
533 are unable to control their basic emotions during straightforward cognitive assessments.  
534 (3) The ontological principle. The ontogeny of social mammals serves as the foundation  
535 for the development of basic emotions, as the source of psychological diseases. (4) The  
536 principle of vulnerability. The susceptibility of individuals to psychiatric diseases varies  
537 according to intrinsically established conditions and adverse circumstances. (5) The  
538 principle of inferential consequences. People pay more attention to an abnormal basic  
539 emotion, engage in introspection to identify their causes. They become skilled at  
540 making inferences about the topic they are pondering, and their inferences can  
541 perpetuate and exacerbate the mental illness.

542 In brief, the Hyper-Emotion theory endorses the notion that individuals  
543 occasionally perform cognitive assessments, which they may consciously recognize,  
544 resulting in an unconscious transition towards a fundamental emotion of heightened  
545 intensity. The episode may be brief or it may intensify into a full-fledged psychological  
546 disease, contingent upon individual constitutional factors and environmental influences.  
547 The theoretical foundation of this study aims to allow adolescents to express their  
548 unconscious emotional perturbations to emotional stimuli from the interactive robot.

### 549 **Data Preprocessing**

550 To ensure that the feature vector dimensions entered into GAME are consistent, we  
551 preprocess the recording data as follows to ensure that the length of the recordings is  
552 the same for all subjects. We set the valid recording duration to 10 seconds as the  
553 average length. If the recording length is longer than the average length, the surplus  
554 frames are truncated, while recordings shorter than the average are zero-padded.  
555  
556

557 Notably, for other data modalities (i.e., inconsistent length of text, face video, and  
558 physiological index), we do not require a preprocessing step due to the inherent  
559 capabilities of the feature extractor in resolving length inconsistencies.

560

### 561 **Single-modal Feature Extraction**

562 The purpose of feature extraction is to retain decent separability (e.g., help GAME  
563 classify data accurately) and reduce computing costs while mapping the sample from a  
564 high-dimensional feature space to a low-dimensional feature space. The followings are  
565 the algorithms used to extract single-modal features or cross-modal features.

#### 566 (1) Feature extraction for audio recordings

567 Mel-scale Frequency Cepstral Coefficients (MFCC)<sup>73</sup> is used as the feature of  
568 acoustic recordings that is commonly used in audio-related tasks like speech  
569 recognition and speaker recognition. An audio is subjected to a rapid Fourier transform,  
570 Mel filter bank, logarithmic operation, discrete offline transform, and dynamic feature  
571 extraction in order to acquire the MFCC feature. We obtain the MFCC feature extracted  
572 by speech-features-module ([https://github.com/jameslyons/python\\_speech\\_features](https://github.com/jameslyons/python_speech_features)),  
573 which is a python package for audio signal processing and audio feature extraction.

574 The calculation of MFCC can be divided into the following steps: first, frame the  
575 signal into brief frames. Under the premise that the audio signal doesn't vary  
576 substantially across small time scales, we confine the signal length into 25 ms, which  
577 is consistent with the acquisition frequency of 16 Khz, corresponding to  $0.025 * 16000 = 400$   
578 frames. We set frame step as 10 ms (160 samples), which allows some  
579 overlap between steps. The first 400 sample frame starts at sample 0, the next 400  
580 sample frame starts at sample 160 etc. until the end of the speech file is reached. The  
581 second step is to calculate the power spectrum of each frame. One set of 12 MFCC  
582 coefficients is retrieved for each frame. Then, the Discrete Fourier Transform (DST)  
583 for each frame will be determined using the following formula:

$$584 \quad S_i(k) = \sum_{n=1}^N s_i(n)h(n)e^{-j2\pi kn/N} \quad 1 \leq k \leq K,$$

585 where  $h(n)$  means the analysis window with  $N$  samples (i.e., hamming window) and  
586  $K$  is the length of the DFT. Additionally,  $s(n)$  means time domain signal, whose  $i$   
587 ranges over the number of frames. The  $S_i(k)$  and  $P_i(k)$  implies the time-domain  
588 frame and the power spectrum of frame  $i$ , respectively. Then, the periodogram-based  
589 power spectral estimate for the speech frame  $s_i(n)$  is given below:

$$590 \quad P_i(k) = \frac{1}{N} |S_i(k)|^2.$$

591 We square the output after taking the complex Fourier transform's absolute value. The  
592 next step is to calculate the Mel-spaced filter bank, take the log for each of the 26 output  
593 from previous step, and finally take DCT of the 26 log filter bank items to obtain 26  
594 cepstral coefficients. Consistent with traditional automatic speech recognition task  
595 settings, we keep the lower 13 of the 26 coefficients as the resulting features.

596 In addition to the conventional speech recognition algorithm for feature extraction ,  
597 we also employ the self-supervised pre-training DL model wav2vec 2.0<sup>74</sup> to embed the  
598 audio. In contrast to other models, wav2vec 2.0 performs the best in many standard  
599 voice tasks<sup>75</sup>. Thus, we employ wav2vec to extract features from audio recordings of

600 adolescents. Wav2vec2.0 encodes speech audio using a multi-layer convolution neural  
601 network and subsequently masks portions of the latent speech representations. The  
602 model is trained using a contrastive manner in which the real latent is differentiated  
603 from fake latent. The latent representations are supplied to a Transformer<sup>76</sup> network to  
604 produce contextualized representations.

#### 605 (2) Feature extraction for textual transcripts

606 For text data, we use **Robustly optimized BERT** approach (RoBERTa)<sup>77</sup> and  
607 PERT<sup>78</sup> to extract the textual feature. These two models yield distinct features due to  
608 differences in architecture and training data from a Chinese corpus. Consequently, we  
609 harness the two models' output as the inputs to improve the robustness and reliability  
610 of GAME in predicting adolescents' mental disorders. RoBERTa and PERT are being  
611 advanced iterations of BERT<sup>79</sup>, exhibiting capability in numerous tasks including text  
612 classification, machine reading comprehension, and text prediction. Based on pre-  
613 trained models, we extract features directly without fine-tuning. RoBERTa is an  
614 improved BERT model that can match or exceed the performance of all post-BERT  
615 methods and it offers a comprehensive evaluation concerning the impact of hyper-  
616 parameter tuning and change of training set size<sup>77</sup>. PERT is a permuted language model  
617 to recover the word orders from a disordered sentence, and the objective of PERT is to  
618 predict the position of the original word, which outperforms other BERT variants on a  
619 few tasks<sup>78</sup>. The amalgamation of PERT and RoBERTa serves to extract the features of  
620 text data from different perspectives.

#### 621 (3) Feature extraction for facial images

622 The features of the face images are extracted using MediaPipe FaceMesh<sup>80</sup>. This  
623 powerful tool, even when presented with single images devoid of depth information, is  
624 capable of furnishing a 3D representation of the human face, comprising 468 points  
625 characterized by 3D coordinates. We use the pre-trained model to generate the features  
626 from each image in the sequence (Supplementary **Fig. 14**), in which the face is resized  
627 to  $256 \times 256$ . The initial processing step entails the application of a facial detector to  
628 delineate a rectangular region encompassing the face, inclusive of vital landmarks such  
629 as eye centers and nose tips. Then the face rectangle is cropped, resized, and fed to a  
630 deep neural network to generate a vector of 3D landmark coordinates.

631 Furthermore, we use MediaPipe Iris<sup>81</sup> to track the eye movements of the volunteer  
632 (Supplementary **Fig. 14**). After MediaPipe FaceMesh detects the face area and eye  
633 landmarks, a DL model is trained to mark subtle positions such as iris position, eye  
634 contour, and pupil location. The position of each eye is represented by a pair of  
635 coordinates. Eye movement can be utilized to infer users' behavior and cognitive status  
636 in human-computer interaction<sup>82</sup>, since pupil response is closely related to cognitive  
637 and emotional processes<sup>83</sup>.

#### 638 (4) Feature extraction for physiological indicators

639 Tsfresh<sup>84</sup> is a Python package for extracting features from time series data, which  
640 employs a repertoire of 63 methods to obtain features, such as absolute energy, the  
641 highest absolute value, etc. The Tsfresh module processes the time series data in three  
642 stages. The first phase is feature extraction, in which the algorithm characterizes the  
643 time series and generates aggregated time series features using the module of feature

644 calculators. Each extracted feature vector is weighted according to their respective p-  
645 values to determine significance in achieving the desired outcome during the feature  
646 significance testing phase. The concluding phase involves a multiple test procedure,  
647 which determines what features need to be retained<sup>85</sup>. The detailed implementation of  
648 feature extraction is described in Supplementary **Methods**.

649

### 650 **Z-Score Normalization**

651 After extracting the modal features from the individual modality data, we transform  
652 them using Z-score normalization to convert the feature vectors into a consistent spatial  
653 dimension. The following formula is used to determine the Z-score in statistics:

654

$$Z = (x - \mu) / \sigma$$

655

656 where, Z means Z-score, x is the original value being evaluated,  $\mu$  denotes the mean  
657 value of all data and  $\sigma$  implies the standard deviation. Cross-modal feature extraction  
658 and multimodal feature fusion are performed after Z-score normalization.

658

### 659 **Cross-modal Feature Extraction**

660 From eight single-modal features standardized by Z-score, we extract cross-modal  
661 features: Relation graph and Attention, in the pursuit of advancing the capabilities of  
662 the GAME. Cross-modal features mine the relationship between various modal features,  
663 assisting GAME to use the correlation among modal features to predict a variety of  
664 mental disorders. The Relation graph is conceptualized as a weighted undirected graph,  
665 wherein each node represents an individual single-modal feature. The weight assigned  
666 to each edge in this graph is determined by the proximity between the respective feature  
667 nodes. Since the length of different unimodal features varies, we apply the Dynamic  
668 Time Warping (DTW)<sup>86</sup> approach to compare the similarity between two time series of  
669 varying lengths or calculate the distance between them. Consequently, the resulting  
670 relation graph is characterized by a vertex set comprising eight nodes and an edge set  
671 comprising 32 weighted edges, all of which are succinctly encapsulated within an  
672  $8 \times 8$  adjacency matrix.

673

674 For the calculation process of DTW, suppose we need to measure the distance  
675 between two example series  $X = \{x_1, x_2, \dots, x_m\}$  and  $Y = \{y_1, y_2, \dots, y_n\}$ . We set  
676  $M(X, Y)$  as the  $m \times n$  point-by-point distance matrix between sequences X and Y,

676

677 where each point (i, j) is distance calculated by  $M_{i,j} = (a_i - b_j)^2$  after the alignment

677

678 between  $x_i$  and  $y_j$  due to length variation. The elements of X and Y are mapped  
679 along a warping path P to minimize the distance between them and P is a group of

679

680 index pairs that make up a matrix traversal, which is defined as:

680

$$P = \langle (e_1, f_1), (e_2, f_2), \dots, (e_s, f_s) \rangle$$

681

682 In order to avoid the problem of combinatorically explosive (i.e., examining every  
683 possible combination), the following prerequisites must be satisfied for a warping path  
684 to be valid: (1) Boundary Condition:  $(e_1, f_1) = (1, 1)$  and  $(e_s, f_s) = (m, n)$ , which  
685 guarantees that the warping path starts at the beginning of both series and terminates at

685

686 the endpoints of them. (2) Monotonicity condition:  $e_i \leq e_{i+1}, 0 < i \leq m$  and  $f_i \leq$   
687  $f_{i+1}, 0 < i \leq n$ , which preserves the chronological sequence of points. (3) Continuity

687 condition:  $e_{i+1} - e_i \leq 1, 0 < i \leq m$  and  $f_{i+1} - f_i \leq 1, 0 < i \leq n$ , which restricts  
688 the forward transitions to nearby points in next time-stage. We define  $\text{dist}(X_{x_i}, Y_{y_i})$  be  
689 the distance between elements at point  $x_i$  of sequence  $X$  and  $y_i$  of sequence  $Y$ . As  
690 a consequence, the distance for optimal path  $P$  is equal to

$$691 \quad D_P(X_{x_i}, Y_{y_i}) = \text{dist}(X_{x_i}, Y_{y_i}) + \min \{D_P(X_{x_{i-1}}, Y_{y_i}), D_P(X_{x_i}, Y_{y_{i-1}}), D_P(X_{x_{i-1}}, Y_{y_{i-1}})\}.$$

692 If we use  $\Theta$  to represent the realm of all potential paths and  $P^*$  is the shortest warping  
693 path. Hence, we can calculate the optimal warping path that

$$694 \quad P^* = \min_{P \in \Theta} (D_P(X, Y)).$$

695 Let  $p_i = M_{X_{e_i}, Y_{f_i}}$  be the distance between elements at position  $e_i$  belong to  $X$  and  
696  $f_i$  of  $Y$ . The DTW distance between two series is obtained by the formula:

$$697 \quad D_{P^*}(X, Y) = \sum_{i=1}^S p_i.$$

698 An exact solution of the best route  $P^*$  can be made using a dynamic programming  
699 approach.

700 With attention mechanism, the model can extract crucial feature, assign each input  
701 component a different weight, and reach more precise judgments. Similarly, we  
702 leverage the DTW method with attention weights, and the detailed process is described  
703 as the following. First, we select one of the single-modal features as the benchmark and  
704 use the DTW technique to determine the distance with the other remaining features. We  
705 use  $d_i$  to denote the distance between any two single-modal features,  $d_i = DTW(M)$ ,  
706  $0 \leq i \leq 7$ , where  $M$  is the feature vector set with eight unimodal features. Second, we  
707 utilize the softmax function convert the distance set  $D = \{d_i\}, 0 \leq i \leq 7$  produced in  
708 the first step into a weight set  $W = \{w_i\}, 0 \leq i \leq 7$  to satisfy the requirements that  
709  $\sum_{i=0}^7 w_i = 1$ . Third, the corresponding feature vector is weighted based on the weight  
710 set obtained in the second stage, and the outcome is then added in bitwise to the  
711 benchmark feature vector. The addition operation is based on the sequence  
712 correspondence in the DTW algorithm, and the dimensionality of the resulting feature  
713 vector is the same as the benchmark. Forth, repeat the same procedures using each of  
714 the eight single-modal features as the reference to generate eight new feature vectors,  
715 and then concatenate them as the attention modal feature.

716

## 717 **Multimodal Feature Fusion and Classification**

### 718 (1) Task-level feature fusion

719 Here we use a simple strategy of averaging all feature vectors including text, audio, and  
720 the face landmarks. The average of eight sentence features is used to describe the  
721 overall features of the text modality, the average of five audio features is used to  
722 describe the features of the audio modality, and the average of multiple face landmarks  
723 is used to represent the face's 3D shape feature. For the iris location in the face image,  
724 we use it directly without any preprocessing before multimodal fusion.

### 725 (2) GAME

726 GAME extracts eight unimodal features from four individual modality data and  
727 creates two novel cross-modal features based on the single-modal features. We then

728 employ EmbraceNet<sup>35</sup> as the backbone network of the multimodal feature fusion  
729 method, and the network structure of GAME is shown in **Figure 2**. EmbraceNet is a  
730 robust multimodal fusion model allowing for excellent compatibility with any network  
731 structure, which considers correlations between various modalities. Additionally,  
732 GAME can handle missing data. There are two main parts in EmbraceNet: the docking  
733 layers and the embracement layer. Docking layers convert the feature vector of a  
734 modality into a format suitable for integration, where the original feature vector is  
735 multiplied with parameter matrix and added by bias matrix. For example, suppose that  
736 there are  $m$  modal features extracted by corresponding network models, the output  
737 vector from the  $k^{th}$  network model will be called  $x^{(k)}$ , where  $1 \leq k \leq m$ . The  $i^{th}$   
738 component of the input vector for the  $k^{th}$  docking layer is written as

$$739 \quad z_i^{(k)} = w_i^{(k)} \cdot x^{(k)} + b_i^{(k)},$$

740 where  $w_i^{(k)}$  and  $b_i^{(k)}$  are weight and bias vector that correspond to the  $k^{th}$  docking  
741 layer, respectively. Finally, the output  $d^{(k)}$  of the  $k^{th}$  docking layer is obtained by  
742 applying an activation function  $f_a$  to  $z_i^{(k)}$ , i.e.,

$$743 \quad d_i^{(k)} = f_a(z_i^{(k)}).$$

744 All the outputs of the docking layers are vectors with  $c$  dimensions, where the  
745 hyper-parameter  $c$  (embracement size) can be configured if necessary (32 in GAME).

746 In the embracement layer, the outputs of the docking layers are fused into a vector  
747 representing all modal information using a probability-based approach as follows.

748 Consider  $r_i = [r_i^{(1)}, r_i^{(2)}, \dots, r_i^{(m)}]^T$ ,  $1 \leq i \leq c$  is a vector obtained from a multinomial  
749 distribution,  $r_i \sim \text{multinomial}(1, p)$ , where  $p = [p_1, p_2, \dots, p_m]$  and  $\sum_{k=1}^m p_k = 1$ .  
750 Only one  $r_i$  equals to 1 in accordance with the definition of the multinomial  
751 distribution, and all other values are equal to 0. The vector  $r^{(k)} = [r_1^{(k)}, r_2^{(k)}, \dots, r_c^{(k)}]^T$   
752 is calculated with the output vector from docking layers  $d^{(k)}$  as

$$753 \quad d'^{(k)} = [d_1'^{(k)}, d_2'^{(k)}, \dots, d_c'^{(k)}]^T = r^{(k)} \circ d^{(k)},$$

754 where  $\circ$  means the Hadamard product, which will multiple the elements in bitwise (i.e.,

755  $d_i'^{(k)} = r_i^{(k)} \cdot d_i^{(k)}$ ). Ultimately, the  $i^{th}$  element of the output vector belonging to the

756 embracement layer  $e = [e_1, e_2, \dots, e_c]^T$  is determined by the following formula:  $e_i =$

757  $\sum_{k=1}^m d_i'^{(k)}$ . The terminal network uses it as an input vector and outputs a final category

758 label for the specified classification task.

759

## 760 **Experimental Evaluation Metrics**

761 In order to comprehensively evaluate the performance of GAME on imbalanced  
762 datasets, we implement a stratified k-fold cross-validation approach, where  $k$  is set as

763 10. Accuracy, weighted F1-score, weighted Precision score, weighted Recall score, and  
764 normalized confusion matrix are calculated. The accuracy can be computed by the  
765 formula:

$$766 \quad Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

767 The F1-score is calculated by Precision score and Recall score. The definitions of the  
768 weighted Precision score and weighted Recall score are listed as the following.

$$769 \quad Precision_i = \frac{TP_i}{TP_i + FP_i}$$

$$770 \quad Precision_{weighted} = \frac{\sum_{i=1}^L (Precision_i \times w_i)}{L}$$

$$771 \quad Recall_i = \frac{TP_i}{TP_i + FN_i}$$

$$772 \quad Recall_{weighted} = \frac{\sum_{i=1}^L (Recall_i \times w_i)}{L}$$

$$773 \quad w_i = \frac{Sn_i}{Tn}$$

774 where  $i$  depicts class index,  $L$  is the total class number,  $TP$  means True positive,  
775  $TN$  is True negative,  $FP$  represents False positive,  $FN$  is False negative,  $Sn$  is  
776 sample number of specific class, and  $Tn$  is the total sample number. The weighted  
777 F1-Score can be determined as

$$778 \quad F1_{weighted} = 2 \times \frac{Precision_{weighted} \times Recall_{weighted}}{Precision_{weighted} + Recall_{weighted}}$$

779 Normalized confusion matrix in cross validation is obtained by averaging each fold of  
780 the confusion matrix and then normalizing the output.

781

## 782 **Acknowledgements**

783 We appreciate the participants in this study for their time and valuable commitment to  
784 this study. We thank Dr. Yongjie Zhou from Shenzhen Mental Health Center for her  
785 time and comments about the labeling of ground truth for each adolescent mental  
786 disorder. This work is supported by funding from the National Natural Science  
787 Foundation of China 31970752; Science, Technology, Innovation Commission of  
788 Shenzhen Municipality JCYJ20190809180003689, JSGG20200225150707332,  
789 JCYJ20220530143014032, ZDSYS20200820165400003,  
790 WDZC20200820173710001, WDZC20200821150704001, JSGG20191129110812708;  
791 Shenzhen Bay Laboratory Open Funding, SZBL2020090501004; Department of  
792 Chemical Engineering-iBHE special cooperation joint fund project, DCE-iBHE-2022-  
793 3; Tsinghua Shenzhen International Graduate School Cross-disciplinary Research and  
794 Innovation Fund Research Plan, JC2022009; and Bureau of Planning, Land and  
795 Resources of Shenzhen Municipality (2022) 207.

796

## 797 **Competing interests**

798 The authors have declared no competing interests or potential conflicts that could have



799 appeared to influence the work reported in this paper.

800

### 801 **Data availability**

802 Due to requirements for ethical approval and the possibility of jeopardizing participant  
803 privacy, we will publish our dataset after feature extraction instead of the original  
804 dataset.

805

### 806 **Code availability**

807 All the code supporting this work will be available at the GitHub repository after  
808 acceptance of manuscript.

809

### 810 **Reference**

- 811 1 Kim-Cohen, J. *et al.* Prior juvenile diagnoses in adults with mental disorder: developmental  
812 follow-back of a prospective-longitudinal cohort. *Archives of general psychiatry* **60**, 709-  
813 717, doi:10.1001/archpsyc.60.7.709 (2003).
- 814 2 Kessler, R. C. *et al.* Lifetime prevalence and age-of-onset distributions of mental disorders  
815 in the World Health Organization's World Mental Health Survey Initiative. *World*  
816 *psychiatry : official journal of the World Psychiatric Association (WPA)* **6**, 168-176 (2007).
- 817 3 Keeley, B. The State of the World's Children 2021: On My Mind--Promoting, Protecting  
818 and Caring for Children's Mental Health. *UNICEF* (2021).
- 819 4 Kieling, C. *et al.* Child and adolescent mental health worldwide: evidence for action. *Lancet*  
820 *(London, England)* **378**, 1515-1525, doi:10.1016/s0140-6736(11)60827-1 (2011).
- 821 5 Delamater, A. M., Guzman, A. & Aparicio, K. Mental health issues in children and  
822 adolescents with chronic illness. *International Journal of Human Rights in Healthcare* **10**,  
823 163-173, doi:10.1108/IJHRH-05-2017-0020 (2017).
- 824 6 Costello, E. J., He, J. P., Sampson, N. A., Kessler, R. C. & Merikangas, K. R. Services for  
825 adolescents with psychiatric disorders: 12-month data from the National Comorbidity  
826 Survey-Adolescent. *Psychiatric services (Washington, D.C.)* **65**, 359-366,  
827 doi:10.1176/appi.ps.201100518 (2014).
- 828 7 Haberer, J. E., Trabin, T. & Klinkman, M. Furthering the reliable and valid measurement of  
829 mental health screening, diagnoses, treatment and outcomes through health information  
830 technology. *General hospital psychiatry* **35**, 349-353,  
831 doi:10.1016/j.genhosppsych.2013.03.009 (2013).
- 832 8 Castiglioni, M. & Laudisa, F. Toward psychiatry as a 'human' science of mind. The case of  
833 depressive disorders in DSM-5. *Frontiers in psychology* **5**, 1517,  
834 doi:10.3389/fpsyg.2014.01517 (2014).
- 835 9 Fakhoury, M. Artificial Intelligence in Psychiatry. *Advances in experimental medicine and*  
836 *biology* **1192**, 119-125, doi:10.1007/978-981-32-9721-0\_6 (2019).
- 837 10 Aguirre Velasco, A., Cruz, I. S. S., Billings, J., Jimenez, M. & Rowe, S. What are the barriers,  
838 facilitators and interventions targeting help-seeking behaviours for common mental  
839 health problems in adolescents? A systematic review. *BMC Psychiatry* **20**, 293,  
840 doi:10.1186/s12888-020-02659-0 (2020).
- 841 11 Rasouli, S., Gupta, G., Ghafurian, M. & Dautenhahn, K. Proposed Applications of Social  
842 Robots in Interventions for Children and Adolescents with Social Anxiety. *Sixteenth*

843 *International Conference on Tangible, Embedded, and Embodied Interaction*, Article 71,  
844 doi:10.1145/3490149.3505575 (2022).

845 12 Abbasi, N. I. *et al.* Can Robots Help in the Evaluation of Mental Wellbeing in Children? An  
846 Empirical Study. *2022 31st IEEE International Conference on Robot and Human Interactive*  
847 *Communication (RO-MAN)*, 1459-1466, doi:10.1109/RO-MAN53752.2022.9900843  
848 (2022).

849 13 Richmond-Rakerd, L. S., D'Souza, S., Milne, B. J., Caspi, A. & Moffitt, T. E. Longitudinal  
850 Associations of Mental Disorders With Physical Diseases and Mortality Among 2.3 Million  
851 New Zealand Citizens. *JAMA network open* **4**, e2033448,  
852 doi:10.1001/jamanetworkopen.2020.33448 (2021).

853 14 McGorry, P. D. *et al.* Designing and scaling up integrated youth mental health care. *World*  
854 *psychiatry : official journal of the World Psychiatric Association (WPA)* **21**, 61-76,  
855 doi:10.1002/wps.20938 (2022).

856 15 Wu, Z. *et al.* Changes of psychotic-like experiences and their association with  
857 anxiety/depression among young adolescents before COVID-19 and after the lockdown  
858 in China. *Schizophrenia Research* **237**, 40-46, doi:10.1016/j.schres.2021.08.020 (2021).

859 16 Wang, J., Li, Y. & He, E. J. P. S. Development and standardization of mental health scale  
860 for middle school students in China. *Psychosoc. Sci* **4**, 15-20 (1997).

861 17 Dong, R.-b. & Dou, K.-y. Changes in physical activity level of adolescents and its  
862 relationship with mental health during regular COVID-19 prevention and control. *Brain*  
863 *and Behavior n/a*, e31116 (2023).

864 18 Desideri, L. *et al.* Using a Humanoid Robot as a Complement to Interventions for Children  
865 with Autism Spectrum Disorder: a Pilot Study. *Advances in Neurodevelopmental Disorders*  
866 **2**, 273-285, doi:10.1007/s41252-018-0066-4 (2018).

867 19 Alves-Oliveira, P. *et al.* Robot-mediated interventions for youth mental health. *Design for*  
868 *Health* **6**, 138-162, doi:10.1080/24735132.2022.2101825 (2022).

869 20 Li, T. W. *et al.* Tell Me About It: Adolescent Self-Disclosure with an Online Robot for Mental  
870 Health. *Companion of the 2023 ACM/IEEE International Conference on Human-Robot*  
871 *Interaction*, 183-187, doi:10.1145/3568294.3580068 (2023).

872 21 Doraiswamy, P. M., Blease, C. & Bodner, K. Artificial intelligence and the future of  
873 psychiatry: Insights from a global physician survey. *Artif Intell Med* **102**, 101753,  
874 doi:10.1016/j.artmed.2019.101753 (2020).

875 22 Dwyer, D. & Koutsouleris, N. Annual Research Review: Translational machine learning for  
876 child and adolescent psychiatry. *Journal of Child Psychology and Psychiatry* **63**, 421-443,  
877 doi:10.1111/jcpp.13545 (2022).

878 23 Moura, I. *et al.* Digital Phenotyping of Mental Health using multimodal sensing of multiple  
879 situations of interest: A Systematic Literature Review. *Journal of Biomedical Informatics*  
880 **138**, 104278, doi:10.1016/j.jbi.2022.104278 (2023).

881 24 Coutts, L. V., Plans, D., Brown, A. W. & Collomosse, J. J. J. o. B. I. Deep learning with  
882 wearable based heart rate variability for prediction of mental and general health. *J Biomed*  
883 *Inform* **112**, 103610, doi:10.1016/j.jbi.2020.103610 (2020).

884 25 Arya, L. & Sethia, D. HRV and GSR as Viable Physiological Markers for Mental Health  
885 Recognition. *2022 14th International Conference on COMMunication Systems &*  
886 *NETworks (COMSNETS)*, 37-42 (2022).

- 887 26 Tiwari, S. & Agarwal, S. J. B. D. A Shrewd Artificial Neural Network-Based Hybrid Model  
888 for Pervasive Stress Detection of Students Using Galvanic Skin Response and  
889 Electrocardiogram Signals. *Big Data* **9**, 427-442, doi:10.1089/big.2020.0256 (2021).
- 890 27 Tariq, Q. *et al.* Mobile detection of autism through machine learning on home video: A  
891 development and prospective validation study. *PLoS medicine* **15**, e1002705,  
892 doi:10.1371/journal.pmed.1002705 (2018).
- 893 28 Cohen, J. *et al.* A feasibility study using a machine learning suicide risk prediction model  
894 based on open-ended interview language in adolescent therapy sessions. *Int J Environ  
895 Res Public Health* **17**, 8187, doi:10.3390/ijerph17218187 (2020).
- 896 29 Saha, K., Yousuf, A., Boyd, R. L., Pennebaker, J. W. & De Choudhury, M. J. S. r. Social media  
897 discussions predict mental health consultations on college campuses. *Sci Rep* **12**, 123,  
898 doi:10.1038/s41598-021-03423-4 (2022).
- 899 30 Huang, Y. *et al.* What Makes Multi-Modal Learning Better than Single (Provably).  
900 *Advances in Neural Information Processing Systems*, 10944-10956,  
901 doi:10.48550/arXiv.2106.04538 (2021).
- 902 31 Tunc, B. *et al.* Diagnostic shifts in autism spectrum disorder can be linked to the fuzzy  
903 nature of the diagnostic boundary: a data-driven approach. *Journal of Child Psychology  
904 and Psychiatry* **62**, 1236-1245, doi:10.1111/jcpp.13406 (2021).
- 905 32 Zhang-James, Y. *et al.* Machine-Learning prediction of comorbid substance use disorders  
906 in ADHD youth using Swedish registry data. *Journal of Child Psychology and Psychiatry*  
907 **61**, 1370-1379, doi:10.1111/jcpp.13226 (2020).
- 908 33 Worthington, M. A. *et al.* Individualized Prediction of Prodromal Symptom Remission for  
909 Youth at Clinical High Risk for Psychosis. *Schizophrenia Bulletin* **48**, 395-404,  
910 doi:10.1093/schbul/sbab115 (2022).
- 911 34 Marti-Puig, P., Capra, C., Vega, D., Llunas, L. & Sole-Casals, J. A Machine Learning  
912 Approach for Predicting Non-Suicidal Self-Injury in Young Adults. *Sensors* **22**,  
913 doi:10.3390/s22134790 (2022).
- 914 35 Choi, J.-H. & Lee, J.-S. EmbraceNet: A robust deep learning architecture for multimodal  
915 classification. *Information Fusion* **51**, 259-270, doi:doi.org/10.1016/j.inffus.2019.02.010  
916 (2019).
- 917 36 Hughes, C. W. & Melson, A. G. in *Handbook of psychiatric measures, 2nd ed.* 251-  
918 308 (American Psychiatric Publishing, Inc., 2008).
- 919 37 Gangemi, A., Tenore, K. & Mancini, F. Two Reasoning Strategies in Patients With  
920 Psychological Illnesses. *Frontiers in psychology* **10**, 2335, doi:10.3389/fpsyg.2019.02335  
921 (2019).
- 922 38 Johnson-Laird, P. N., Mancini, F. & Gangemi, A. A hyper-emotion theory of psychological  
923 illnesses. *Psychological review* **113**, 822-841, doi:10.1037/0033-295x.113.4.822 (2006).
- 924 39 Schumann, G. *et al.* The IMAGEN study: reinforcement-related behaviour in normal brain  
925 function and psychopathology. *Molecular Psychiatry* **15**, 1128-1139,  
926 doi:10.1038/mp.2010.4 (2010).
- 927 40 Karcher, N. R. & Barch, D. M. The ABCD study: understanding the development of risk for  
928 mental and physical health outcomes. *Neuropsychopharmacology* **46**, 131-142,  
929 doi:10.1038/s41386-020-0736-6 (2021).
- 930 41 Rezapour, M. & Hansen, L. A machine learning analysis of COVID-19 mental health data.

931 *Scientific Reports* **12**, 14965, doi:10.1038/s41598-022-19314-1 (2022).

932 42 Mohamed, E. S. *et al.* A hybrid mental health prediction model using Support Vector  
933 Machine, Multilayer Perceptron, and Random Forest algorithms. *Healthcare Analytics* **3**,  
934 100185 (2023).

935 43 Garriga, R. *et al.* Machine learning model to predict mental health crises from electronic  
936 health records. *Nature Medicine* **28**, 1240-1248, doi:10.1038/s41591-022-01811-5  
937 (2022).

938 44 Montgomery, J., Hendry, J., Wilson, J. A., Deary, I. J. & MacKenzie, K. Pragmatic detection  
939 of anxiety and depression in a prospective cohort of voice outpatient clinic attenders.  
940 *Clinical Otolaryngology* **41**, 2-7, doi:10.1111/coa.12459 (2016).

941 45 Mo, L., Li, H. & Zhu, T. Exploring the Suicide Mechanism Path of High-Suicide-Risk  
942 Adolescents-Based on Weibo Text Analysis. *Int J Environ Res Public Health* **19**,  
943 doi:10.3390/ijerph191811495 (2022).

944 46 Bhardwaj, V. *et al.* Machine Learning for Endometrial Cancer Prediction and  
945 Prognostication. *Frontiers in oncology* **12**, 852746, doi:10.3389/fonc.2022.852746 (2022).

946 47 Xie, Y. *et al.* Stroke prediction from electrocardiograms by deep neural network.  
947 *Multimedia Tools and Applications* **80**, 17291-17297, doi:10.1007/s11042-020-10043-z  
948 (2021).

949 48 Githinji, B. *et al.* Multidimensional Hypergraph on Delineated Retinal Features for  
950 Pathological Myopia Task. *25th International Conference on Medical Image Computing  
951 and Computer Assisted Intervention (MICCAI)* **13432**, 550-559, doi:10.1007/978-3-031-  
952 16434-7\_53 (2022).

953 49 Yang, J. *et al.* Prediction of HER2-positive breast cancer recurrence and metastasis risk  
954 from histopathological images and clinical information via multimodal deep learning.  
955 *Computational and Structural Biotechnology Journal* **20**, 333-342,  
956 doi:10.1016/j.csbj.2021.12.028 (2022).

957 50 Arbabshirani, M. R., Plis, S., Sui, J. & Calhoun, V. D. Single subject prediction of brain  
958 disorders in neuroimaging: Promises and pitfalls. *Neuroimage* **145**, 137-165,  
959 doi:10.1016/j.neuroimage.2016.02.079 (2017).

960 51 Ashwini, B. Robot Assisted Diagnosis of Autism in Children. *Proceedings of the 2020  
961 International Conference on Multimodal Interaction*, 728-732,  
962 doi:10.1145/3382507.3421162 (2020).

963 52 Uban, A.-S., Chulvi, B. & Rosso, P. Explainability of Depression Detection on Social Media:  
964 From Deep Learning Models to Psychological Interpretations and Multimodality. *Early  
965 Detection of Mental Health Disorders by Social Media Monitoring: The First Five Years of  
966 the eRisk Project*, 289-320, doi:10.1007/978-3-031-04431-1\_13 (2022).

967 53 Stone, L. L., Otten, R., Engels, R. C. M. E., Vermulst, A. A. & Janssens, J. M. A. M.  
968 Psychometric Properties of the Parent and Teacher Versions of the Strengths and  
969 Difficulties Questionnaire for 4- to 12-Year-Olds: A Review. *Clinical child and family  
970 psychology review* **13**, 254-274, doi:10.1007/s10567-010-0071-2 (2010).

971 54 Castellanos-Ryan, N., O'Leary-Barrett, M., Sully, L. & Conrod, P. Sensitivity and Specificity  
972 of a Brief Personality Screening Instrument in Predicting Future Substance Use, Emotional,  
973 and Behavioral Problems: 18-Month Predictive Validity of the Substance Use Risk Profile  
974 Scale. *Alcoholism: Clinical and Experimental Research* **37**, E281-E290,

975 doi:doi.org/10.1111/j.1530-0277.2012.01931.x (2013).

976 55 Gili, M. *et al.* Mental disorders as risk factors for suicidal behavior in young people: A  
977 meta-analysis and systematic review of longitudinal studies. *Journal of Affective Disorders*  
978 **245**, 152-162, doi:doi.org/10.1016/j.jad.2018.10.115 (2019).

979 56 Hirschfeld, R. M. The Comorbidity of Major Depression and Anxiety Disorders: Recognition  
980 and Management in Primary Care. *Primary care companion to the Journal of clinical*  
981 *psychiatry* **3**, 244-254, doi:10.4088/pcc.v03n0609 (2001).

982 57 Davids, A. & Parenti, A. N. Time orientation and interpersonal relations of emotionally  
983 disturbed and normal children. *The Journal of Abnormal and Social Psychology* **57**, 299-  
984 305, doi:10.1037/h0047687 (1958).

985 58 Chaddad, A., Peng, J., Xu, J. & Bouridane, A. Survey of Explainable AI Techniques in  
986 Healthcare. *Sensors (Basel, Switzerland)* **23**, doi:10.3390/s23020634 (2023).

987 59 Duijndam, S., Karreman, A., Denollet, J. & Kupper, N. Physiological and emotional  
988 responses to evaluative stress in socially inhibited young adults. *Biological Psychology* **149**,  
989 doi:10.1016/j.biopsycho.2019.107811 (2020).

990 60 Byck, G. R. *et al.* Effect of housing relocation and neighborhood environment on  
991 adolescent mental and behavioral health. *Journal of Child Psychology and Psychiatry* **56**,  
992 1185-1193, doi:10.1111/jcpp.12386 (2015).

993 61 Gross, J. J. & Levenson, R. W. EMOTION ELICITATION USING FILMS. *Cognition & Emotion*  
994 **9**, 87-108, doi:10.1080/02699939508408966 (1995).

995 62 Schaefer, A., Nils, F., Sanchez, X. & Philippot, P. Assessing the effectiveness of a large  
996 database of emotion-eliciting films: A new tool for emotion researchers. *Cognition &*  
997 *Emotion* **24**, 1153-1172, doi:10.1080/02699930903274322 (2010).

998 63 Ma, M. *et al.* Smil: Multimodal learning with severely missing modality. *Proceedings of the*  
999 *AAAI Conference on Artificial Intelligence* **35**, 2302-2310 (2021).

1000 64 Dolci, G. *et al.* in *2023 IEEE International Conference on Acoustics, Speech, and Signal*  
1001 *Processing Workshops (ICASSPW)*. 1-5.

1002 65 Hillege, R. H. L., Lo, J. C., Janssen, C. P. & Romeijn, N. The Mental Machine: Classifying  
1003 Mental Workload State from Unobtrusive Heart Rate-Measures Using Machine Learning.  
1004 *Adaptive Instructional Systems. Second International Conference, AIS 2020. Held as Part*  
1005 *of the 22nd HCI International Conference, HCII 2020. Proceedings. Lecture Notes in*  
1006 *Computer Science (LNCS 12214)*, 330-349, doi:10.1007/978-3-030-50788-6\_24 (2020).

1007 66 Davila, M. I., Lewis, G. F. & Porges, S. W. The PhysioCam: A Novel Non-Contact Sensor to  
1008 Measure Heart Rate Variability in Clinical and Field Applications. *Frontiers in Public Health*  
1009 **5**, doi:10.3389/fpubh.2017.00300 (2017).

1010 67 Nishikawa, M. *et al.* in *43rd Annual International Conference of the IEEE-Engineering-in-*  
1011 *Medicine-and-Biology-Society (IEEE EMBC)*. 7016-7019 (2021).

1012 68 Camras, L. EMOTION - A PSYCHOEVOLUTIONARY SYNTHESIS - PLUTCHIK,R. *American*  
1013 *Journal of Psychology* **93**, 751-753, doi:10.2307/1422394 (1980).

1014 69 Zhu, H. B., Han, G. J., Shu, L. & Zhao, H. ArvaNet: Deep Recurrent Architecture for PPG-  
1015 Based Negative Mental-State Monitoring. *Ieee Transactions on Computational Social*  
1016 *Systems* **8**, 179-190, doi:10.1109/tcss.2020.2977715 (2021).

1017 70 Bonicatto, S., Dew, M. A., Soria, J. J. & Seghezso, M. E. Validity and reliability of Symptom  
1018 Checklist '90 (SCL90) in an Argentine population sample. *Social Psychiatry and Psychiatric*

1019            *Epidemiology* **32**, 332-338, doi:10.1007/bf00805438 (1997).

1020    71    Luo, Y. *et al.* Mental Health Problems and Associated Factors in Chinese High School  
1021            Students in Henan Province: A Cross-Sectional Study. *Int J Environ Res Public Health* **17**,  
1022            doi:10.3390/ijerph17165944 (2020).

1023    72    Chen, Y. *et al.* Social Identity, Core Self-Evaluation, School Adaptation, and Mental Health  
1024            Problems in Migrant Children in China: A Chain Mediation Model. *Int J Environ Res Public*  
1025            *Health* **19**, doi:10.3390/ijerph192416645 (2022).

1026    73    das, A., Jena, M. R. & Barik, K. K. Mel-Frequency Cepstral Coefficient (MFCC) - a Novel  
1027            Method for Speaker Recognition. *Digital Technologies* **1**, 1-3, doi:10.12691/dt-1-1-1  
1028            (2014).

1029    74    Baevski, A., Zhou, H., Mohamed, A. & Auli, M. wav2vec 2.0: a framework for self-  
1030            supervised learning of speech representations. *Proceedings of the 34th International*  
1031            *Conference on Neural Information Processing Systems*, Article 1044 (2020).

1032    75    Hsu, W.-N. *et al.* Robust wav2vec 2.0: Analyzing domain shift in self-supervised pre-  
1033            training. *arXiv preprint*, doi:arXiv:2104.01027 (2021).

1034    76    Vaswani, A. *et al.* Attention is all you need. *Proceedings of the 31st International*  
1035            *Conference on Neural Information Processing Systems*, 6000-6010 (2017).

1036    77    Yinhan, L. *et al.* RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint*,  
1037            13 pp.-13 pp., doi:arXiv:1907.11692 (2019).

1038    78    Cui, Y., Yang, Z. & Liu, T. PERT: Pre-training BERT with Permuted Language Model. *arXiv*  
1039            *preprint*, doi:arXiv:2203.06906 (2022).

1040    79    Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional  
1041            Transformers for Language Understanding. *arXiv preprint*, 4171-4186,  
1042            doi:arXiv:1810.04805 (2018).

1043    80    Kartynnik, Y., Ablavatski, A., Grishchenko, I. & Grundmann, M. Real-time facial surface  
1044            geometry from monocular video on mobile GPUs. *arXiv preprint*,  
1045            doi:arXiv:1907.06724 (2019).

1046    81    Ablavatski, A., Vakunov, A., Grishchenko, I., Raveendran, K. & Zhdanovich, M. J. a. p. a.  
1047            Real-time Pupil Tracking from Monocular Video for Digital Puppetry. *arXiv preprint*,  
1048            doi:arXiv:2006.11341 (2020).

1049    82    Zheng, W.-L., Liu, W., Lu, Y., Lu, B.-L. & Cichocki, A. EmotionMeter: A Multimodal  
1050            Framework for Recognizing Human Emotions. *Ieee Transactions on Cybernetics* **49**, 1110-  
1051            1122, doi:10.1109/tcyb.2018.2797176 (2019).

1052    83    Zekveld, A. A., Heslenfeld, D. J., Johnsrude, I. S., Versfeld, N. J. & Kramer, S. E. The eye as  
1053            a window to the listening brain: Neural correlates of pupil size as a measure of cognitive  
1054            listening load. *Neuroimage* **101**, 76-86, doi:10.1016/j.neuroimage.2014.06.069 (2014).

1055    84    Christ, M., Braun, N., Neuffer, J. & Kempa-Liehr, A. W. Time Series FeatuRe Extraction on  
1056            basis of Scalable Hypothesis tests (tsfresh – A Python package). **307**, 72-77,  
1057            doi:10.1016/j.neucom.2018.03.067 (2018).

1058    85    Christ, M., Kempa-Liehr, A. W. & Feindt, M. J. a. p. a. Distributed and parallel time series  
1059            feature extraction for industrial big data applications. *Asian Machine Learning Conference*  
1060            *(ACML) 2016, Workshop on Learning on Big Data (WLBD), Hamilton (New Zealand)* (2016).

1061    86    Lines, J. & Bagnall, A. Time series classification with ensembles of elastic distance measures.  
1062            *Data Mining and Knowledge Discovery* **29**, 565-592, doi:10.1007/s10618-014-0361-2

1063  
1064

(2015).

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplement.docx](#)