

DepGuard: Depression Recognition and Episode Monitoring System With a Ubiquitous Wrist-Worn Device

Yufei Zhang¹, Shuo Jin, Wenting Kuang, Yuda Zheng, Qifeng Song, Changhe Fan²,
Yongpan Zou³, *Member, IEEE*, Victor C. M. Leung⁴, *Life Fellow, IEEE*, and Kaishun Wu⁵, *Fellow, IEEE*

Abstract—Depression significantly impacts mental health, severely disrupting patients’ daily lives. During depressive episodes, individuals may experience symptoms such as excessive guilt, self-harm, and suicidal ideation. Compared to proprietary devices like brain electrode caps, wearable technologies for depression detection have gained attention due to their affordability and portability—enabling real-time monitoring of depressive states. However, challenges such as low-quality data from ubiquitous devices, individual variability, and the complexity of multimodal physiological signal analysis limit model generalizability. To address these issues, we present *DepGuard*, a novel ubiquitous wearable system for depression assessment based on multimodal physiological signals. *DepGuard* performs a two-stage detection process: depression recognition and real-time episode monitoring. For depression recognition, we propose an unsupervised domain adaptation method to reduce the domain gap between source and target subjects. For episode monitoring, we employ a few-shot learning strategy to enable personalized modeling. Both approaches enhance cross-subject generalization. Our system achieves 90.75% accuracy in cross-subject depression recognition using 30 unlabeled

samples per target subject, and 93.52% accuracy in episode monitoring using 15 labeled samples per class.

Index Terms—Wearable devices, depression recognition, unsupervised domain adaptation, depressive episode monitoring.

I. INTRODUCTION

DEPRESSION is a common mental disorder affecting approximately 280 million people worldwide [1]. It is characterized by persistent sadness, insomnia, loss of interest, and reduced ability to carry out daily activities for at least two weeks [2]. Without timely treatment, individuals are at increased risk of suicidal behavior. The Lancet Psychiatry, a leading journal in the field, emphasizes that early recognition and intervention targeting lifestyle habits and risk factors are critical for effective prevention [3]. However, continuous supervision by doctors or family members is not feasible, making early detection of depressive episodes essential to improving patient outcomes and reducing the risk of self-harm or suicide. Real-time monitoring, coupled with timely communication to healthcare professionals, plays a vital role in ensuring prompt and appropriate care. This proactive approach not only helps prevent adverse events but also enhances overall mental health management.

In recent years, depression recognition research has increasingly combined individual-level insights with advanced technologies to enable objective, passive assessment of depressive states. Specifically, given the characteristic patterns of negative affect and social withdrawal in individuals with depression, mental states can be inferred through the analysis of lifestyle data, including sleep patterns [4], [5], [6], social connectivity [7], [8], [9], behavioral records [10], [11], [12], [13]. However, behavior-based methods for depression detection are often hindered by individuals intentionally or unintentionally concealing their symptoms. Similarly, traditional self-report scales suffer from subjectivity and typically lack robust quantitative analysis. Physiological signal monitoring offers a promising alternative. Depression-related symptoms such as emotional instability and physical fatigue often manifest as measurable physiological changes that are difficult to consciously control, making them less susceptible to self-report bias [14]. High-end commercial wearable devices can capture rich physiological

Received 30 December 2024; revised 18 June 2025; accepted 16 July 2025. Date of publication 21 July 2025; date of current version 3 December 2025. This work was supported in part by China NSFC under Grant 62172286 and Grant U2001207, in part by Key Project of Department of Education of Guangdong Province under Grant 2024ZDZX1011, in part by Guangdong Basic and Applied Basic Research Foundation under Grant 2022A1515011509, in part by the Guangdong Provincial Key Lab of Integrated Communication, Sensing and Computation for Ubiquitous Internet of Things under Grant 2023B1212010007, in part by the Project of DEGP under Grant 2023KCXTD042, in part by Guangdong “Pearl River Talent Recruitment Program” under Grant 2019ZT08X603, in part by Guangdong “Pearl River Talent Plan” under Grant 2019JC01X235, in part by 111 Center under Grant D25008, and in part by Shenzhen Science and Technology Foundation under Grant ZDSYS20190902092853047. Recommended for acceptance by W. Wang. (Corresponding author: Yongpan Zou.)

This work involved human subjects in its research. Approval of all ethical and experimental procedures and protocols was granted by Institutional Review Board (IRB), Faculty of Medicine, Shenzhen University under Application No. 202500165, and performed in line with the Declaration of Helsinki.

Yufei Zhang, Shuo Jin, Wenting Kuang, Yuda Zheng, Qifeng Song, and Yongpan Zou are with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen 518060, China (e-mail: yfyaya@foxmail.com; jinshuo2022@email.szu.edu.cn; kuangwenting2023@email.szu.edu.cn; song qifeng2021@email.szu.edu.cn; 2021150056@email.szu.edu.cn; yongpan@szy.edu.cn).

Changhe Fan is with the Department of Psychology, Guangdong Second Provincial General Hospital, Guangzhou 510317, China (e-mail: changhefan@yahoo.com).

Victor C. M. Leung is with the Artificial Intelligence Research Institute, Shenzhen MSU-BIT University, Shenzhen 518172, China (e-mail: vleung@ieee.org).

Kaishun Wu is with the Information Hub, Hong Kong University of Science and Technology, Guangzhou 510317, China (e-mail: wuks@hkust-gz.edu.cn).

Digital Object Identifier 10.1109/TMC.2025.3591096

data—such as electroencephalography (EEG), electrocardiography (ECG), and magnetoencephalography (MEG) [15], but their high cost, limited accessibility, and operational complexity restrict widespread application. In contrast, ubiquitous wearable devices offer low-cost and portable solutions, but they tend to capture only shallow physiological patterns, which weakens their effectiveness in depression recognition and episode monitoring. Consequently, proprietary equipment remains the standard for high-accuracy detection [16], [17], [18], despite its impracticality for everyday use [19], [20].

Recent research has begun to explore the potential of ubiquitous devices for mental health monitoring [21], [22], [23], [24], though their application in depression detection remains in its early stages. For instance, Hassantabar et al. [21] introduce MHDeep, a framework that integrates wearable sensor data with artificial neural networks to classify mental health disorders. Their model utilizes multimodal inputs such as galvanic skin response, skin temperature, inter-beat interval, and tri-axial acceleration, collected from smartwatches and smartphones. Similarly, Hu et al. [23] develop an ensemble machine learning model that classifies users as depressed or non-depressed based on sleep data from wearable devices. Ahmed et al. [25] propose a multimodal framework that combines convolutional neural networks, attention mechanisms, and random forests to assess affective states and depression. Despite these advancements, several key challenges remain unresolved.

C1: How can depressive episodes be detected quickly to provide timely interventions? Much of the existing research focuses solely on depression recognition at the model level, using either professional or consumer-grade equipment, with limited attention to real-time monitoring of depressive episodes. However, the progression from diagnosis to ongoing monitoring is a continuous process, inherently linking detection with treatment. These tasks are causally and practically connected in real-world settings. To alleviate the substantial burden placed on healthcare providers and family members during inpatient care, a comprehensive systems approach is needed—one that enables timely intervention and mitigates the psychological risks associated with depressive episodes.

C2: How can we improve model adaptability to achieve high-accuracy detection for new subjects? Existing research often overlooks cross-domain challenges, limiting model performance on unseen individuals. Due to the uniqueness of individual physiological traits and emotional responses, even those with similar mental health conditions may exhibit different physiological signal patterns. This challenge is compounded by limited data and scarce annotations, highlighting the need for models to learn more discriminative and generalizable patterns.

To address the aforementioned challenges, we introduce *DepGuard*, a depression assessment system that analyzes multimodal physiological signals to identify both depressive disorders and ongoing depressive episodes. To tackle the issue of real-time episode monitoring, *DepGuard* integrates a dedicated module that operates in the system's second phase. Users wear a wrist-worn device that continuously collects physiological data, which is then transmitted via Bluetooth to a paired mobile application. A lightweight and personalized model is deployed within the

app to perform real-time predictions and deliver results to the user. Given the device's limited computational resources, we employ a knowledge distillation strategy to compress the model, thereby reducing its complexity while maintaining predictive performance.

To address the second challenge, we design task-specific strategies that align with the unique requirements of each objective. For depression recognition, we utilize an unsupervised domain adaptation approach to transfer knowledge from the source domain (i.e., labeled data from known subjects) to the target domain (i.e., unseen subject), thereby improving generalization. For depressive episode monitoring, we apply few-shot learning to enhance adaptability with minimal labeled data. By integrating transfer learning and few-shot learning, our system effectively mitigates inter-individual variability and enables robust, cross-subject depression assessment and monitoring.

Our main contributions can be summarized as follows.

- We present *DepGuard*, a wearable depression assessment system designed for both depression diagnosis and real-time monitoring of depressive episodes. Evaluated under cross-subject testing on a real-world dataset, *DepGuard* achieves 90.75% accuracy in depression recognition and 93.52% in depressive episode monitoring.
- We propose an unsupervised domain adaptation method for depression recognition that bridges the gap between source and target domains. This approach minimizes data and training overhead for new subjects while maintaining high recognition accuracy.
- We propose a few-shot learning approach combined with a knowledge distillation strategy for depressive episode monitoring, aiming to build high-performance personalized models on edge devices with minimal computational cost.
- To facilitate this line of the research,¹ we collect a real-world dataset from individuals diagnosed with depression by certified psychiatrists at our partner hospital.

The remainder of the paper is conducted as follows. Section II reviews the related works about different types of devices and cross-subject detection of depression. Subsequently, Section III introduces the two-stage procedure of *DepGuard*. Then, Section IV conducts the implementation of our proposals. Section V evaluates the effectiveness of *DepGuard*. Finally, Sections VI and VII discuss and conclude our works.

II. RELATED WORKS

A. Depression Detection With Different Devices

Proprietary devices such as EEG [16], [17], [26], ECG [18], [27], MEG [28], functional near-infrared spectroscopy (fNIRS) [29] and functional magnetic resonance imaging (fMRI) [30] are commonly used to detect physiological signals. The high resolution, precision, and stability of these devices allow for the detection of subtle physiological changes, providing researchers with reliable, high-quality data. However,

¹Our source code is available at <https://github.com/ISAC-GROUP/AFFC-DepGuard>.

these devices require complex operations when worn. For example, Huang et al. [28] use a bulky whole-head 306-channel Elekta Neuronmag Vectorview to capture MEG signals. Pange et al. [31] using EEG and ECG signals require conductive gel and power connections, making the setup inconvenient for users.

For depression detection based on portable devices, several studies [11] have indicated that the great significance of utilizing sensor data from smartphones and wearable devices, such as Fitbit and Microsoft Band, can effectively represent the features of depression defined in the DSM-5 diagnostic criteria [2]. These sensors data encompass various dimensions including phone usage patterns, acceleration, GPS, sleep duration, and audio recordings, providing a comprehensive understanding of subjects physical activity, social interactions, and attention levels in daily life. Analyzing this data opens new possibilities for precision medicine in diagnosing and treating depression. As a result, portable devices are explored for physiological signal detection and further utilized to predict depression. Specifically, Tazawa et al. [32] apply machine learning to screen and assess depression severity using Silmee W20 wristband data, achieving 76% accuracy. Moshe et al. [33] utilize Oura Ring data on activity, sleep, and heart rate variability to predict depression, revealing significant correlations with depressive symptoms. Moreover, Hassantabar et al. [21] collect physiological signals from both the E4 smartwatch and Samsung smartphone, achieving an average accuracy of 87.3% in distinguishing healthy and depressive instances. Based on a lightweight TN1012/ST Pulse Transducer, the research [22] develops a pulse rate variability detection model to identify depression with the highest accuracy of 98.46%. However, data collected by portable devices face technical constraints, including low temporal resolution and noise susceptibility, which collectively impact prediction reliability. These limitations motivate the development of a wrist-worn device equipped with multi-source physiological sensors. Furthermore, due to inter-individual variability, research on depression recognition remains underdeveloped and requires further progress.

B. Cross-Subject Depression Detection

Individual-independent detection tasks have become a central research challenge. In depression detection, physiological and psychological variations, such as differences in neural activity, mood fluctuations, and cognitive responses, pose substantial obstacles to the generalization of models trained on source domain data when applied to target domains. As a result, to enhance the accuracy and robustness of cross-subject tasks, some of the research is directed towards employing advanced techniques such as domain adaptation [34], [35] and domain generalization [36], which aim to mitigate the domain shift between source and target data distributions, improving the accuracy and robustness of cross-subject models. For instance, Zhang et al. [36] propose a novel EEG-based Graph Neural Network for depression detection. In their proposal, domain generalization based on adversarial training is adapted to the model, and a secondary subject partitioning method is proposed to group subjects with similar data distributions into the same domain with a shared domain label. Fang et al. [37] present an unsupervised cross-domain fMRI

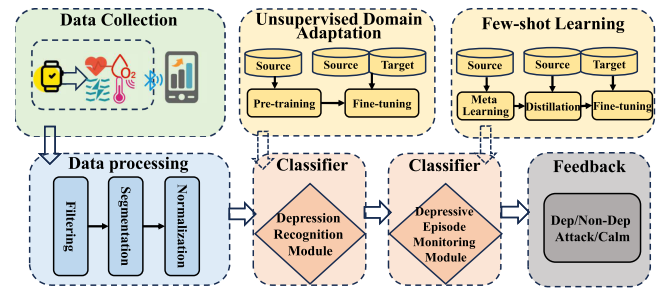


Fig. 1. The overview of the data-processing pipeline.

adaptation framework (UFA-Net), addressing heterogeneity in major depressive disorder detection. This approach leverages graph convolutional networks and maximum mean discrepancy for feature alignment, improving model generalizability across domains without target labels. Chen et al. [35] introduce a Graph Neural Network-based Semi-supervised Domain Adaptation (GNN-SDA) framework for major depressive disorder detection. The GNN-SDA addresses semantic misalignment and class imbalance issues by aligning domains at an individual level and optimizing pseudo-labels, enhancing the model's ability to generalize across different datasets without relying on target labels. Zhang et al. [38] introduce a multi-classification model, DepL-GCN, that employs a graph convolutional network (GCN) and domain generalization to classify depression levels from EEG data. This model addresses subject variability and sample imbalance across categories by integrating a penalty coefficient for smaller classes, enhancing its ability to generalize and accurately detect varying degrees of depression without relying on target domain labels.

As for our design, we adopt targeted cross-domain techniques tailored to the distinct characteristics of the two task modules: depression recognition and depressive episode monitoring. The first phase focuses on depression recognition. Given the absence of labeled data in the target domain, we employ an unsupervised domain adaptation module, UDA-DR, which transfers knowledge from the source domain while effectively minimizing the distribution gap without relying on target domain labels. In the second phase, a small amount of labeled data is available due to the subjects diagnosed with depression. To maximize the utility of this data, we propose the FSL-DEM module based on few-shot learning, which enables rapid adaptation to new individuals and supports personalized monitoring of depressive episodes. By combining unsupervised and few-shot learning, *DepGuard* enhances model generalization and adapts to varying levels of data availability. This dual approach is a key advantage of the *DepGuard* system in addressing cross-individual depression detection challenges.

III. SYSTEM DESIGN

In this section, we introduce *DepGuard* system for recognizing depression and monitoring depressive episodes. As shown in Fig. 1, the system pipeline begins with the collection of multimodal physiological signal data via a wrist-worn device, which is transmitted to software for specialized preprocessing.

The processed data is then used for initial-stage depression detection. If depression is detected, the system proceeds to monitor depressive episodes in the second stage. Otherwise, the results are directly displayed. The depressive episode monitoring task involves real-time analysis of the multimodal physiological signals collected from the patient, providing an assessment of whether the patient is currently in a depressive state.

A. Data Preprocessing

Due to the prohibitive costs of off-the-shelf commercial devices for collecting physiological signals and the lack of available data interfaces, we develop our hardware and software solutions to facilitate data collection and other tasks.

Heart rate signal: To mitigate PPG signal drift caused by varying ambient light, we apply a Detrended Fluctuation Analysis (DFA) with a time window length of 20 to remove trend components. The PPG signal comprises two main frequency bands: low-frequency components (0.05 to 0.5 Hz), which reflect sympathetic nerve activity, fluid regulation, and changes in respiration and blood pressure; and high-frequency components (0.5 to 4 Hz), which are primarily associated with parasympathetic nerve activity and short-term heart rate variability (HRV). To remove noise and other interferences, we apply a Butterworth band-pass filter with a frequency range of 0.05 to 4 Hz.

Blood oxygen saturation signal: On the collected infrared and red light signal data, we first apply a Short-time Fourier Transform (STFT) to analyze the signal's frequency components. Then, we use a Butterworth band-pass filter (0.1 ~ 4 Hz) to remove noise and improve signal quality.

Galvanic skin response signal: GSR signals primarily show low-frequency conductivity changes. To focus on physiological arousal and sweat gland activity, we apply STFT and a second-order Butterworth filter with a cutoff frequency set at 0.3 Hz.

Skin temperature signal: The trend in skin temperature signal (SKT) is relatively subtle. We apply an STFT and find that the primary components are predominantly within the effective frequency range. Additionally, power frequency interference can be introduced into the sensor signal, affecting the accuracy of the skin temperature sensor. To mitigate power frequency interference, we use a Butterworth filter to truncate the frequency of the skin temperature signal.

After processing the physiological signals, we use sliding processing with a window size of 1 and a stride size of 1. And then, (1) is used to normalize each data channel to a range of -1 to 1 .

$$y_{norm} = 2 \frac{x - x_{min}}{x_{max} - x_{min}} - 1 \quad (1)$$

B. Unsupervised Domain Adaptation for Depression Recognition

Poor performance in cross-subject depression recognition is primarily attributed to the lack of data from new subjects during the training phase, which fails to capture the heterogeneity of individual physiological signals. This diversity is crucial for accurate model training. Without incorporating this range of individual differences, the model struggles to adapt to the unique

physiological patterns of new subjects, thus reducing its effectiveness in recognizing depression across individuals. Additionally, depression is a complex mental illness typically diagnosed through clinical evaluations, making the acquisition of labeled data challenging. The nature of depression recognition prevents us from acquiring labeled target-domain (i.e., target subject) data in the fine-tuning stage, leaving us to depend solely on unlabeled data. Therefore, we propose an unsupervised domain adaptation for depression recognition (UDA-DR) to enhance cross-subject performance. Fig. 2 shows the structure of UDA-DR, which adopts a two-phase training process of pre-training followed by fine-tuning. The method primarily consists of three modules: encoder, class classifier, and domain classifier. The encoder ensures that the class classifier performs well during pre-training and that the domain classifier cannot distinguish between source and target domain data during fine-tuning. The class classifier aims to correctly categorize the data based on features extracted by the encoder, while the domain classifier distinguishes whether the data originates from the source or target domain.

1) *Model Building:* By processing data from five channels, we extract time-domain feature samples. Fig. 2 presents a multi-channel temporal network based on LSTM, CNN and attention mechanism (LCMTN), serving as the backbone network for Sections III-B and III-C. The LCMTN mainly consists of an encoder and classifier, where the classifier uses a Multilayer Perceptron (MLP). We first apply a Long Short-Term Memory (LSTM) network for the feature encoder to extract local temporal features from each channel.

Subsequently, to obtain the coupling of temporal dependencies, we deploy a multi-head self-attention mechanism [39] in LCMTN to enhance the derived representation of temporal features. The multi-head self-attention mechanism enables the model to simultaneously attend to different segments of the input sequence, capturing diverse temporal dependencies and enhancing the global representation of physiological signals over time. Concretely, as shown in (2), given a feature representation $S \in \mathcal{R}^{L \times C}$ fed into linear layers to generate query space, key space and value space, represents learnable weighted matrices, L and C denote the time series and the number of channels, respectively. Next, one head of attention output is formulated as shown in (3), where the softmax function is applied to the scaled dot product of the query and key matrices to produce attention scores. These scores are then used to weigh the values in the value matrix $V(S)$, resulting in an attention head output that captures the importance of different temporal components across the sequence. In (4), the output of multi-head self-attention $MHA(S)$ is a repeated process for multiple heads to jointly enable the model to focus on various aspects of the input sequence, where $W^O \in \mathcal{R}^{C \times C}$ denotes a learnable weighted matrix, $Concat$ is a concatenated operation, and k is the number of attention heads.

$$Q(S) = S \cdot W_Q, K(S) = S \cdot W_K, V(S) = S \cdot W_V \quad (2)$$

$$h_j = \text{Softmax} \left(\frac{Q(S)K^T(S)}{\sqrt{d_K}} \right) V(S) \quad (3)$$

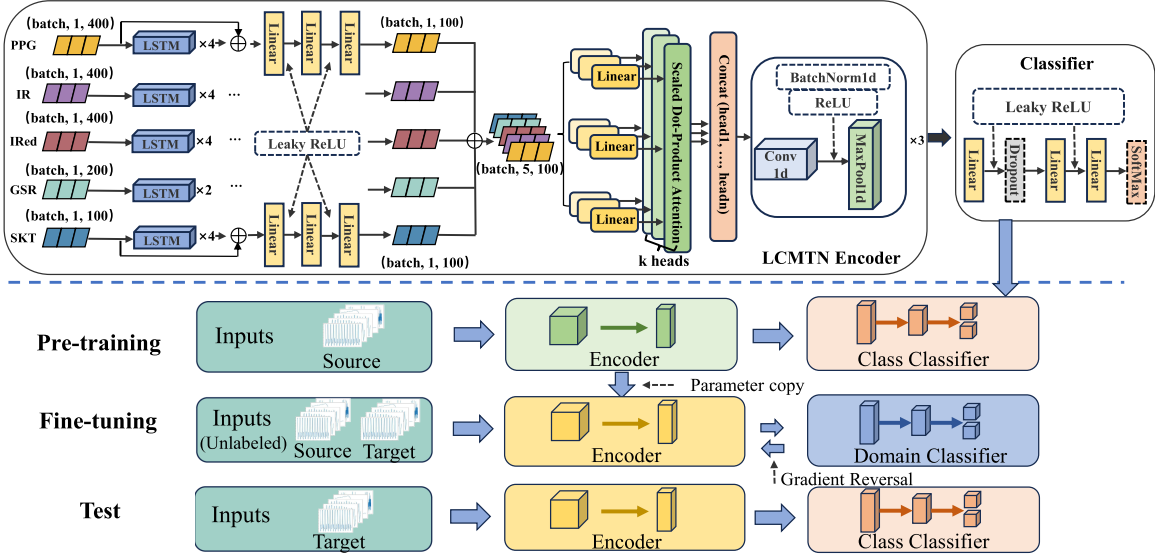


Fig. 2. The architecture of UDA-DR network.

$$\text{MHA}(Q(S), K(S), V(S)) = \text{Concat}(h_1, h_2, \dots, h_k) W^O \quad (4)$$

Finally, we need to integrate features from various channels effectively, ensuring they complement each other for more comprehensive feature information. Thus, we employ a Convolutional Neural Network to precisely capture the correlations between different channels, providing a richer feature representation for the fusion process.

2) *Pre-Training Phase*: The primary objective of this phase is to train LCMTN on large-scale datasets, aiming to obtain a high-performance encoder and class classifier for the task of depression recognition. In the pre-training phase, we only utilize the data of source domain subjects for supervised training. During the training of the network, the data of source domain subjects is partitioned into a training set, validation set and test set. The model parameters are updated through backpropagation using the cross-entropy loss.

3) *Fine-Tuning Phase*: This phase is crucial for enhancing cross-subject performance in depression recognition. During Phase I, we train a high-performance encoder and class classifier optimized for depression recognition. The class classifier demonstrates strong classification capabilities for this task. However, the encoder is not exposed to the data features of the target subject, resulting in poor performance when used directly by the target subject. To address this, we employ unsupervised domain adaptation, utilizing unlabeled data from both source and target domain subjects to fine-tune the encoder. By confusing source and target domain features, the encoder learns domain-invariant features, enabling the class classifier to perform more effectively on target subject data.

Specifically, we mark the source domain data as one category of labels and the target domain data as another category of labels before starting fine-tuning. Then, we introduce a domain classifier using an MLP and integrate it with the encoder obtained in phase I to form an integrated network. The primary task of

this network is to determine whether the input data originates from the source or target domain. A Gradient Reversal Layer (GRL) is placed between the encoder and domain classifier to align the encoder with the target domain data. The GRL inverts the gradient flow during backpropagation, effectively confusing the distinction between source and target domains. This design mitigates domain differences, improving the encoder's adaptability to the target subject data. During the forward pass, the network applies an identity transformation ((5)), while in the backward pass, it reverses the gradient direction ((6)). Here, I denotes the identity matrix, $-I$ indicates reversed gradient flow and λ is a hyperparameter that controls the gradient magnitude during backpropagation.

$$R_\lambda(x) = x \quad (5)$$

$$\frac{dR_\lambda(x)}{dx} = -I \quad (6)$$

At the beginning of the fine-tuning, we input both the source domain and target domain data into the network at a certain ratio in each iteration. When passing through the GRL, it undergoes an identity transformation as usual. The loss function of the entire network is then formulated as:

$$\text{Loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{i,c} \log(p_{i,c}) - \frac{1}{M} \sum_{j=1}^M \sum_{d=1}^D z_{j,d} \log(q_{j,d}) \quad (7)$$

where the first term $-\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{i,c} \log(p_{i,c})$ denotes domain classifier loss L_{source} and the second term $\frac{1}{M} \sum_{j=1}^M \sum_{d=1}^D z_{j,d} \log(q_{j,d})$ denotes the target domain loss L_{target} . N and M denote the number of source domain samples and target domain samples, respectively. C represents the number of classes (depression or non-depression). D represents the number of domains (source domain or target domain). $y_{i,c}$ and $p_{i,c}$ are the true domain label indicator and predicted probability of sample i belonging to class c , respectively. $z_{j,d}$ and $q_{j,d}$

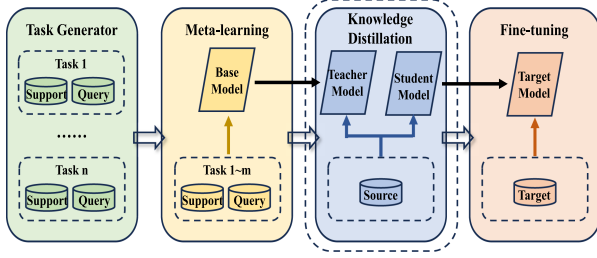


Fig. 3. The architecture of FSL-DEM.

are the true domain label indicator and predicted probability of the sample j belonging to domain d , respectively. During backpropagation, the GRL causes the parameter gradients of the encoder to be reversed. The gradient update formulas for the domain classifier and the encoder are as follows:

$$\theta_y = \theta_y - \mu \frac{\partial L_y}{\partial \theta_y} \quad (8)$$

$$\theta_f = \theta_f - \left(-\mu \frac{\partial L_y}{\partial \theta_f} \right) \quad (9)$$

Here, θ_y represents the parameters of the domain classifier and θ_f represents the parameters of the encoder. Eqs. (8) and (9) show that θ_y updates normally, but θ_f gradients are reversed.

4) *Testing Phase*: After pre-training and fine-tuning phases, the classifier and the encoder are obtained. We combine these to form a personalized depression recognition model for a new subject. This approach not only adapts well to individual differences but also facilitates the creation of personalized models that accurately fit new subjects with minimal data, thereby promoting personalized depression detection.

C. Few-Shot Learning for Depressive Episode Monitoring

Depressive episode monitoring also presents a key challenge: individual differences among subjects reduce cross-subject performance. To address this, strategies such as unsupervised domain adaptation and few-shot learning are commonly used. Although unsupervised domain adaptation is flexible for handling distribution differences between source and target domains, it may be less effective than few-shot learning methods in some cases. In contrast, by leveraging the limited labeled samples in the target domain, few-shot learning can directly guide the model to learn specific patterns in the target domain. Meanwhile, the monitoring task differs from depression recognition in that it can yield limited labeled target domain data for enhancing cross-domain performance. To this end, we propose a few-shot learning method for depressive episode monitoring (FSL-DEM), as shown in Fig. 3. The FSL-DEM consists of four modules: task generator, meta-learning, knowledge distillation, and fine-tuning.

1) *Task Generator*: Given that this task employs the meta-learning method in few-shot learning, it is worth noting that the basic unit of meta-learning training is the meta-task. The primary goal of the task generator is to produce representative and well-defined meta-tasks, ensuring excellent model performance

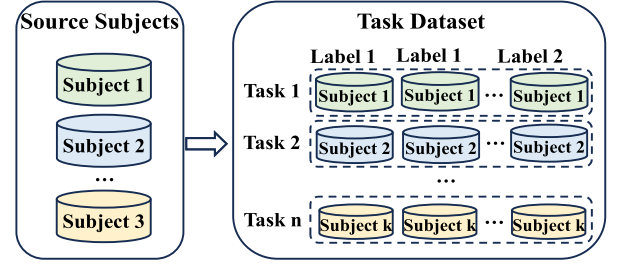


Fig. 4. Subject-based task generator.

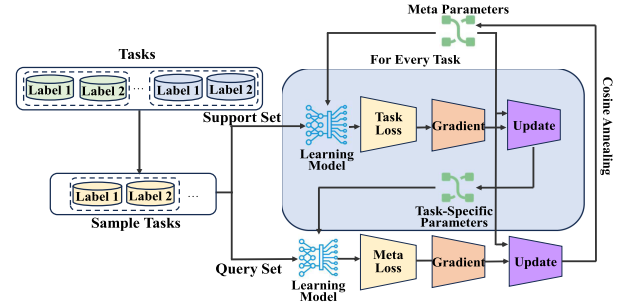


Fig. 5. The Meta-learning training process.

during training. A common strategy for generating meta-tasks involves random selection. This approach aggregates data from all source domain subjects, shuffles it, and then selects a subset to form each meta-task. It is important to note that although we use data from source domain subjects in meta-learning, there are still significant differences in the data of the same categories among different subjects. In fact, the data from each individual can be considered as a unique classification task dataset. This means that in the source domain data, each individual has its unique characteristic and pattern, forming multiple independent classification tasks. As shown in Fig. 4, we treat each source domain subject as an independent dataset and randomly select a certain number of data samples for each category to form a meta-task. In this manner, each source domain subject generates multiple meta-tasks, with an equal number of samples for each category within each meta-task. Meanwhile, we further divide each meta-task into a support set and a query set. The support set is used to train a temporary model for the corresponding meta-task, while the query set is used with the temporary model to compute the loss for updating the base model.

2) *Meta Learning*: The goal of this phase is to obtain well-initialized parameters so that the model can quickly adjust its parameters when facing new subjects, enabling it to adapt more quickly and effectively. Each meta-task consists of a support set and a query set, of which the support set comprises only a small portion. Additionally, the number of samples for each category remains consistent within both the support set and the query set. Fig. 5 illustrates the meta-learning training process, which can be divided into two stages: the inner training stage and the outer training stage.

For the inner training stage, we extract a subset of meta-tasks as an inner training batch and only use the support sets of

sampled tasks for model training. At the start of training, each meta-task independently copies the base model to use as its temporary model, allowing the temporary model to fit some extent on the support set of each task. Subsequently, each meta-task updates its temporary model parameters using the cross-entropy loss function. Each meta-task has an independent loss calculation process that does not interfere with each other, and its parameter gradient update formula is:

$$\theta'_i = \theta - \alpha \nabla_{\theta} L_{T_i}(f_{\theta}) \quad (10)$$

where θ represents the model parameters, α is the learning rate during inner training, $L_{T_i}(f_{\theta})$ is the loss function on task T_i , and $\nabla_{\theta} L_{T_i}(f_{\theta})$ is the gradient of the loss function with respect to the model parameters. In the gradient update, second-order derivatives are generally involved. However, the calculation of the second-order derivatives increases the computational complexity and cost. Therefore, to reduce computational cost, the second-order partial derivative in the formula is approximated to zero, and only the first-order derivative is calculated in practical applications.

The outer training stage is primarily evaluated using a query set of specific tasks and updates the base model parameters. In the inner training stage, we obtain the updated temporary model parameters corresponding to each meta-task. We use the query set of each meta-task for forward propagation in the corresponding temporary model, obtaining the loss for each meta-task. As shown in (11), we sum up these loss to obtain a total loss, which is used to update the parameters of the base model. Subsequently, cosine annealing is employed to continuously adjust the learning rate, which is used to update the parameters of the base model during the outer training stage. Since the goal of meta-learning is to find well-initialized parameters that enable rapid adaptation to new tasks, it is essential to dynamically adjust the learning rate during training to ensure the model converges quickly on different tasks. Cosine annealing precisely offers this dynamic adjustment mechanism, helping the model better escape local optima during the learning process and improving generalization ability. As shown in (12), cosine annealing reduces the learning rate by using a cosine function.

$$L_{total} = \sum_{i=1}^m L_{T_i}(f_{\theta}) \quad (11)$$

$$\beta_t = \beta_{\min} + \frac{1}{2} (\beta_{\max} - \beta_{\min}) \left(1 + \cos \left(\frac{T_{cur}}{T_{\max}} \pi \right) \right) \quad (12)$$

where β_t is the learning rate at step t , β_{\max} and β_{\min} are the maximum and minimum values of the learning rate, defining the range of the learning rate. T_{cur} refers to the current iteration count, and T_{\max} refers to the total number of iterations within one cycle. This formula is designed to make the learning rate decrease slowly at the beginning of a cycle, then drop rapidly, and finally decrease slowly again.

3) *Knowledge Distillation*: This phase of knowledge distillation serves as an optional module. The task of depressive episode monitoring requires real-time performance. If this monitoring module is embedded in devices with limited hardware resources,

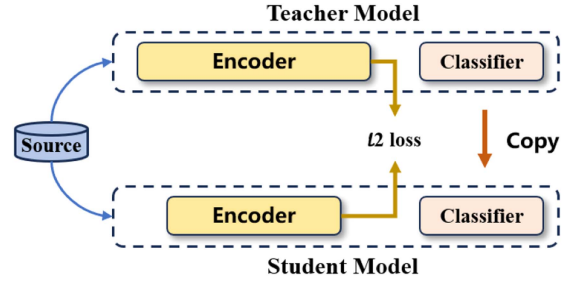


Fig. 6. Feature-based knowledge distillation.

the knowledge distillation module becomes particularly important. Generally, there are three methods of knowledge distillation: feature-based knowledge distillation, response-based knowledge distillation, and relation-based knowledge distillation. When knowledge distillation targets the encoder, the smaller model inherits the representational capacity of the larger model while avoiding overfitting. In contrast, distilling the classifier may cause the student model to overfit to the teacher's predictive behavior, potentially disregarding the true data distribution. To address this, we adopt a feature-based knowledge distillation strategy. As illustrated in Fig. 6, we use the base model trained with meta-learning as the teacher, and a redesigned, lightweight encoder with reduced complexity and fewer parameters as the student. The input data still comes from the source domain, and the loss between the feature vectors output by the encoder of the teacher model and the encoder of the student model is calculated using (13). Subsequently, only the parameters of the encoder in the student model are updated using the obtained loss. The encoder of the trained student model is combined with the classifier of the teacher model to form a new student model.

$$Loss = \frac{1}{n} \sum_{i=1}^n (f_{teacher}(x_i) - f_{student}(x_i))^2 \quad (13)$$

where $f_{teacher}(x_i)$ is the feature vector output by the encoder of teacher model, $f_{student}(x_i)$ is the feature vector output by the encoder of student model, and n is the length of the feature vectors. The formula calculates the ℓ_2 loss between $f_{teacher}(x_i)$ and $f_{student}(x_i)$. When updating the encoder parameters of the student model through this method, the focus is on minimizing the difference in feature representations.

In the knowledge distillation module, we introduce a parameter τ to control the degree of distillation. This parameter determines the extent to which knowledge from the teacher model is transferred to the student model. A higher τ signifies that the student model prioritizes the integration of knowledge distilled from the teacher model, whereas a lower value conversely leads the student model to rely more heavily on the original training dataset for parameter updates. The parameter τ is systematically optimized in Section V-F3 to balance robust generalization capability of the student model with mitigation of over-reliance on the teacher's supervisory signals.

4) *Fine-Tuning Phase*: After meta-learning or knowledge distillation training, the base model has acquired well-initialized

TABLE I
SPECIFICATIONS OF DIFFERENT SENSORS

Signal type	Sensor model	Sampling rate	Data type	Price
Heart rate	Pulse sensor	400 Hz	Raw analog data	3.425 \$
Blood oxygen saturation	Max30102	400 Hz	Red light and infrared light intensity data	0.959 \$
Galvanic skin response	Grove GSR	200 Hz	Voltage simulation data	8.219 \$
Skin temperature	LMT70	100 Hz	Voltage simulation data	3.425 \$

parameters. In this phase, the new subject provides a small amount of data to fine-tune the base model, further adapting it to the characteristics of the target domain data. We divide the fine-tuning process into two stages based on the number of iterations. In the first stage, we update the model parameters with a high learning rate, but the model has not yet reached its best convergence position. In the second stage, a small learning rate is used to update the model parameters to further optimize the model towards the optimal solution and take into account the stability of optimization. In this way, each new subject obtains a customized model tailored to individual characteristics, thereby better adapting to their personalized needs in the depressive episode monitoring task.

IV. SYSTEM IMPLEMENTATION

A. Hardware Design

The hardware part of the *DepGuard* is a wrist-worn device, including four parts: control module, sensor module, communication module, and power module. The control module utilizes STM32, which offers precise data processing and maintains stable performance in low-power mode. The sensor module integrates four ubiquitous biometric sensors for collecting data such as heart rate, blood oxygen saturation, galvanic skin response, and skin temperature.

Table I presents the sensors and their parameters to ensure efficient and accurate acquisition of various environmental data by the *DepGuard*. The communication module uses the Bluetooth BLE5.0 version protocol. The power module employs the SY8089DC-DC buck converter and a TP4054 charging circuit for a maximum 500 mA charge, effectively managing power for the wristband device.

Fig. 7 presents the overall hardware architecture of *DepGuard*. The system comprises a custom-designed motherboard integrated with four ubiquitous biometric sensors, forming a compact, wrist-worn device. The onboard microprocessor collects physiological signals and transmits the data to the companion app via Bluetooth.

B. Software Development

We develop a mobile app for smartphones to complement the *DepGuard* hardware, offering features like physiological monitoring, self-assessment, and depression tracking, making up the full *DepGuard* system. Fig. 8(a)–(e) shows the main interface of the app, highlighting its three core functions. Fig. 8(a) shows the main interface for monitoring physiological data with Bluetooth

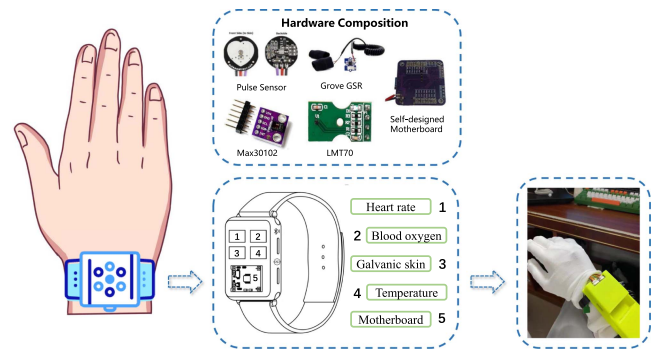


Fig. 7. The design of our wrist-worn device.

and real-time display. The app connects to the hardware via Bluetooth, receives real-time data and stores it for later analysis. When the subjects want to know more about specific physiological data, they can click the icon on the interface to access the corresponding display page. Fig. 8(b) and (c) visualize the heart rate and blood oxygen saturation, plotting their trends and displaying live data. Fig. 8(d) shows the self-assessment interface with personal info and questionnaires including PHQ-9, GAD-7, BDI-II and so on. Users can complete these questionnaires, and their scores, along with severity information for some scales, will be shown on their personal page. Fig. 8(e) is the depressive episode recording interface, which aims to collect the time periods when individuals with depression experience the symptoms described in their medical records. It is important to note that the functionality of this page is only used when we are collecting data on depressive episodes of patients.

The aforementioned interface is integrated into a mobile app and tested on the Huawei Mate 40 Pro smartphone, powered by HiSilicon Kirin 9000 chipset, 4400 mAh battery, and 8 GB RAM. Integrating software and hardware enables *DepGuard* to systematically record the subject's physiological and psychological conditions.

V. EXPERIMENT AND EVALUATION

A. Data Collection

We have undertaken a 3-month data collection and pre-processing process to facilitate the execution of this study. The experimental participants consist of two distinct groups: individuals without depression and those with depression. Before the experiment, all participants agree to sign the Institutional Review Board (IRB) agreement [40]. The data collection process for these two groups diverges; For the group without depression, potential participants are recruited from Shenzhen University. Each participant fills out PHQ-9, GAD-7, BDI-II, and other questionnaires on the software before the experiment. Afterward, participants are selected based on questionnaire scores and criteria to confirm they are “healthy.” Then, we choose a quiet conference room for data collection from non-depressed individuals, keeping them free from distractions. Before starting, participants should get used to the *DepGuard* system to avoid any impact on their physical or mental condition. Subsequently,

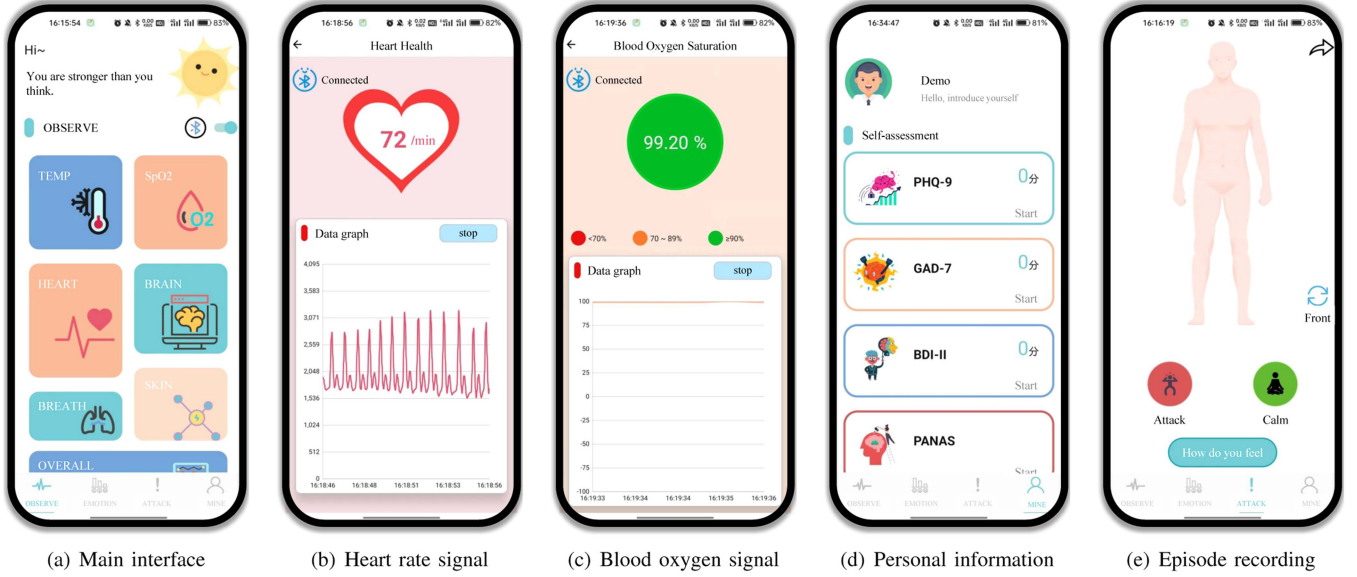


Fig. 8. The main interfaces of our mobile application.

they are given a 1-minute buffer to regulate their emotions to ensure reaching a calm state. On the other side, we collect physiological data from clinically diagnosed depressed patients in Guangdong Second Provincial General Hospital to ensure data authenticity. Patients are informed about the study and sign consent forms. They learn about the *DepGuard* system and are briefed on collection precautions.

Depressive episode data are collected only for patients with diagnosed depression. In the formal phase of data collection, we collect physiological data when the patient is in a state of both physical and mental tranquility. When the patient experiences an episode, they click the ‘Attack’ button on the corresponding interface of the software. When the patient returns to a calm state from the episode, they click the ‘Calm’ button. In this way, the segment of the depressive episode is recorded. Notably, the process of collecting data on depressive episodes described above is used only for the training phase of the model. In real scenarios, the *DepGuard* system collects physiological signals in real-time without clicking on the ‘Attack’ or ‘Calm’ button. Besides, during all of the above data collection procedures, we ask the subjects to remain in a resting state in order to prevent body movement artifacts from being generated. Data on depressive episodes is especially collected in the presence of a healthcare professional. Finally, we recruit 35 participants (20 males, 15 females) with 17 individuals from the healthy control group and 18 patients diagnosed with depression. The ages of the participants range from 14 to 53 years old, with the mean and standard deviation of 23.14 and 7.06, respectively. As for their occupations, twenty-eight participants are still in education, four are currently workers, technicians, or self-employed, and one is retired. The remaining participant does not provide occupational information.

The length of data for subjects with depression recognition and depressive episode detection is inconsistent, ranging from 594 to 7444 seconds. Concretely, the total amount of data for

depression recognition is 42206 seconds. And the total amount of data for the depressive episodes is 38659 seconds, of which 48.7% (18839 seconds) are depressive samples.

B. Evaluation Metrics

For depression recognition, we assess performance using accuracy, detection rate, and misdiagnosis rate. Accuracy, as represented in Equation 14, is the ratio of correctly predicted samples to the total number of samples. Detection rate, as represented in Equation 15, is the ratio of correctly predicted depression samples to the total number of depression samples. Misdiagnosis rate, as represented in Equation 16, is the ratio of incorrectly predicted non-depression samples to the total number of non-depression samples.

$$Accuracy = \frac{n_{pre}}{n_{true}} \cdot 100\% \quad (14)$$

$$Detection\ rate = \frac{n_{pre=dep}}{n_{true=dep}} \cdot 100\% \quad (15)$$

$$Misdiagnosis\ rate = \left(1 - \frac{n_{pre=non-dep}}{n_{true=non-dep}}\right) \cdot 100\% \quad (16)$$

For depressive episode monitoring, we use precision, recall, F1 score, and accuracy as performance evaluation metrics. Among them, precision, recall, and F1 score are all calculated using the macro-average method, which means taking the average of evaluation metrics for each category.

C. Details of Evaluation

In both tasks, the single-subject setting involves randomly partitioning each subject’s data into training (60%), validation (20%), and testing (20%) sets. For cross-subject evaluation, we employ Leave-One-Subject-Out Cross-Validation (LOSO-CV), a specific form of K-Fold cross-validation where K equals the

number of subjects ($K = 35$ for depression recognition; $K = 17$ for depressive episode monitoring). In each iteration, one subject is designated as the test set, while the remaining $K - 1$ subjects constitute the training set. This exhaustive approach ensures that each subject is used once as the test set, providing a comprehensive assessment of the model's generalization across individuals.

All evaluations are implemented with the environment of Intel(R) Xeon(R) CPU E5-2686 v4 @ 2.30 GHz, NVIDIA Tesla A6000 GPU. We harness the computational prowess of the cuDNN library to expedite GPU-accelerated processing. The training process involves the Adam optimizer [41], with a learning rate of $1e - 4$, a batch size of 128, and 200 epochs to ensure robust convergence. The attention head is set to 4 in the MHA module. Random dropping ratio of 0.3 in the Dropout layer.

D. Baseline Methods

We compare the performance of *DepGuard* with the following baseline methods.

- *DANN* [42] is a classic unsupervised domain adaptation method based on adversarial learning. Its main idea is to reduce the distribution discrepancy between the source domain and the target domain by introducing domain adversarial training, thereby improving the model's generalization ability on the target domain.
- *ADDA* [43] is an unsupervised domain adaptation method that combines GAN loss. It employs a different training strategy with DANN to reduce the distribution discrepancy between the source domain and the target domain, thus improving the model's generalization performance.
- *DFA-MCD* [44] is a state-of-the-art unsupervised domain adaptation method that introduces a Gaussian-guided latent alignment approach to align the latent feature distributions of the two domains under the guidance of the prior distribution.
- *SWL-Adapt* [45] is an unsupervised domain adaptation method tailored for temporal input features, incorporating sample weight learning. It calculates sample weights according to the classification loss and domain discrimination loss of each sample with a parameterized network.
- *CDFSL-V* [46] is an effective few-shot learning method for recognizing new categories with only a few labeled examples. The method employs a masked autoencoder-based self-supervised training objective to learn from both source and target data in a self-supervised manner. Then, a progressive curriculum balances learning the discriminative information from the source dataset with the generic information learned from the target domain.
- *PN* [47] is a state-of-the-art few-shot learning method for the image dataset CUB-200-2011, where a classifier must generalize to new classes not seen in the training set, given only a small number of examples of each new class. It learns a metric space in which classification can be performed by computing distances to prototype representations of each class.

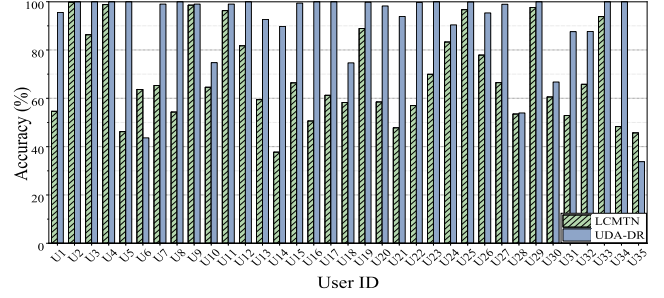


Fig. 9. The cross-subject results of depression subjects.

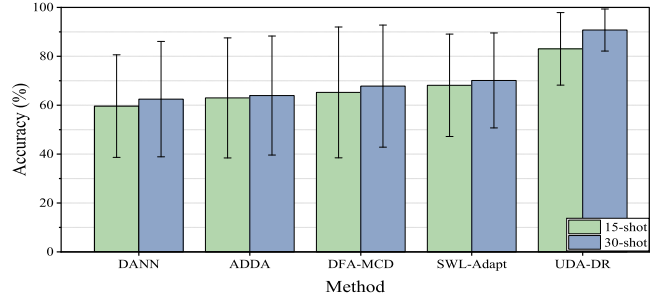


Fig. 10. The results of cross-subject for different methods.

- *PN(-T)* is an ablation method for the PN based on a random task generation strategy.
- *FD(-A)* is an ablation method for the FSL-DEM based on using only easily obtainable non-episode samples.
- *FD(-F)* is an ablation method for the FSL-DEM based on not using fine-tuning module and employing meta-learning training.
- *FD(-T)* is an ablation method for the FSL-DEM based on random task generation strategy.

E. Depression Recognition Results

1) *Performance of Cross-Subject*: We evaluate cross-subject depression recognition using LCMTN and UDA-DR, as depicted in Fig. 9. where U1 to U17 denote non-depressed subjects, and U18 - U35 denote depressed subjects. The overall average accuracy of LCMTN is 68.85%, while that of UDA-DR is 90.75%. UDA-DR enhances cross-subject recognition accuracy by 21.9% over LCMTN. Notably, the accuracy for the non-depression subject U6 and the depression subject U35 decreased. For U6, the depression scale scores are near the threshold for diagnosis. And for U35, the professional psychologist diagnoses fewer and milder symptoms. Moreover, Fig. 10 compares the cross-subject recognition performance of DANN [42], ADDA [43], DFA-MCD [44], SWL-Adapt [45], and UDA-DR under various shot conditions. UDA-DR excels in the 15-shot and 30-shot target domain scenarios due to its effectiveness with limited unlabeled data.

2) *Analysis of Performance Influencing Factors*: Furthermore, we test four key factors for optimizing the cross-subject recognition performance of UDA-DR: the number of target subject data samples and the number of fine-tuning iterations,

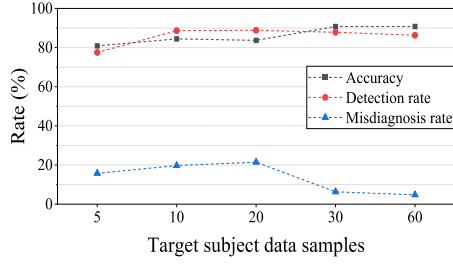


Fig. 11. The performance of different target subject data samples.

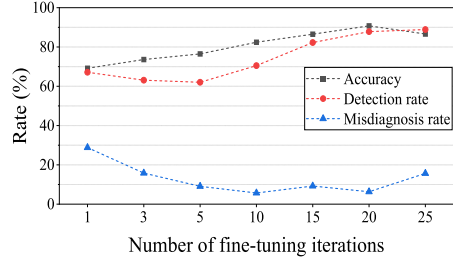


Fig. 12. The performance with the number of fine-tuning iterations.

the number of source domain subjects, and the batch proportion size between the source domain and the target domain.

The number of target subject data samples: Fig. 11 displays the metrics performance for target subject unlabeled data across 5, 10, 20, 30, and 60 shots. Accuracy increases with the increase in the amount of unlabeled data provided to the target subjects. With 5 shots, the accuracy is 80.89%, and the accuracy remains stable and unchanged after 30 shots. Although the detection rate at 30 and 60 shots is slightly lower than at 20 shots, the misdiagnosis rate is significantly better than at 20 shots. Thus, we set a minimum of 30 shots of unlabeled data per target subject.

The number of fine-tuning iterations: In the fine-tuning phase, we fine-tune the target subject data samples of 30 shots with 1, 3, 5, 10, 15, 20, and 25 iterations respectively. Fig. 12 shows the accuracy, detection rate, and misdiagnosis rate results for different fine-tuning numbers. With the continuous increase in the number of fine-tuning, the accuracy shows a gradual upward trend. When fine-tuning is performed 20 times, the accuracy reaches its optimal state. However, further increasing the number of fine-tuning led to a downward trend in accuracy performance. This is because the increase in the number of fine-tuning causes the model to overlearn the specific local structure in the new subject data while ignoring the overall pattern of the data. Although the detection rate of fine-tuning 20 times is slightly lower than that of fine-tuning 25 times, the misdiagnosis rate performance is significantly better than fine-tuning 25 times. Therefore, we fixed the number of fine-tunings to 20 times.

The number of source-domain subjects: The diversity of subjects is closely related to the generalization performance of the model. In the pre-training phase, apart from each individual subject, we randomly select 3, 5, 7, 10, 15 subjects respectively from the remaining subjects in each category as source domain subjects. As shown in Fig. 14, as the number of source domain

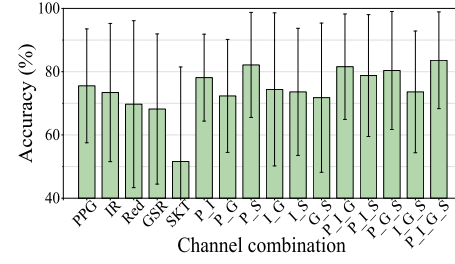


Fig. 13. The results of different channel combinations.

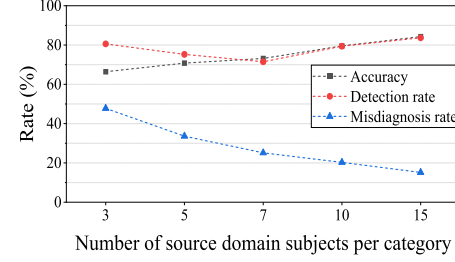


Fig. 14. The performance of different numbers of source domain subjects per category.

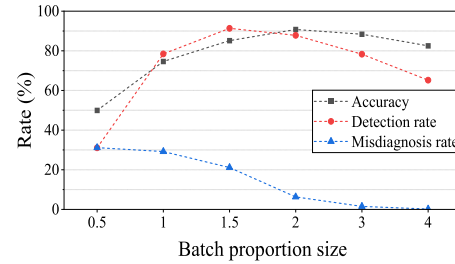


Fig. 15. The performance of different batch proportion between source and target domain.

subjects is continuously increased, the accuracy is gradually improving, and the misdiagnosis rate keeps dropping. When three source domain subjects are provided for each category to all source domain subjects except the target subject, the recognition accuracy increases from 66.4% to 90.75%, and its performance gradually nears the single-subject performance under the LCMTN. It further indicates that although subject diversity is related to model generalization performance, the number of source domain subjects currently used in UDA-DR is already sufficient. Even if the number of source domain subjects is further increased, the space for improvement in recognition performance is limited while the data cost continues to increase.

The batch proportion across domains: In domain adaptation, the batch proportion refers to the ratio of source to target domain samples during fine-tuning, balancing the model's ability to learn discriminative features from the source domain with its adaptability to the target domain. In this phase, we set the batch size of the target domain to a fixed 20, and then respectively set the batch proportion size of the source domain to the target domain as 0.5, 1, 1.5, 2, 3, and 4. As shown in Fig. 15, as the batch

TABLE III
THE PERFORMANCE OF LCMTN FOR SINGLE-SUBJECT AND CROSS-SUBJECT

	Single-subject	Cross-subject
Precision	99.28%	56.39%
Recall	99.28%	56.16%
F1 score	99.27%	53.93%
Accuracy	99.28%	55.78%

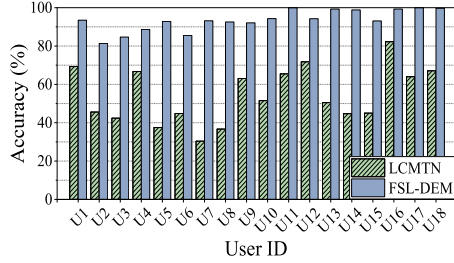


Fig. 17. The results of cross-subject for LCMTN and FSL-DEM.

single-modality signals, limiting its ability to comprehensively capture an individual's peripheral physiological state. In the work [52], manually extracted features are used to construct dynamic representations, followed by a random forest model for prediction. Similarly, the work [48] employs a basic deep neural network, which struggles to effectively fuse multimodal physiological signals. However, their high cost, limited accessibility, and operational complexity significantly hinder real-world application. In contrast, our approach offers superior portability, broader applicability, and enhanced performance, making it more suitable for practical deployment.

F. Depressive Episode Monitoring Results

1) *Performance of Single-Subject and Cross-Subject:* In order to verify the feasibility of depressive episode monitoring, we only conduct single-subject and cross-subject depressive episode monitoring evaluations through the LCMTN. As shown in Table III, we observe that all metrics for single-subject results are vastly higher than cross-subject results, with an accuracy difference of 43.5%. The accuracy range for single-subject results is around 99.28%, indicating a consistently high performance. For cross-subject results, the range is from 53.93% to 55.78%. These results indicate the feasibility of depressive episode monitoring. However, we need to solve the issue of subject variability to improve cross-subject performance for depressive episode monitoring.

To evaluate the cross-subject performance of FSL-DEM, we use the LCMTN and the FSL-DEM method designed in Section III-C. Fig. 17 shows the accuracy of LCMTN and FSL-DEM for each subject. It can be clearly seen that the performance of FSL-DEM is better than that of LCMTN across all subjects. The average accuracy of FSL-DEM is 93.52%, representing a 39.14% increase compared to LCMTN. Among them, subject U7 shows the largest improvement, with accuracy increasing from 30.39% to 93.17%.

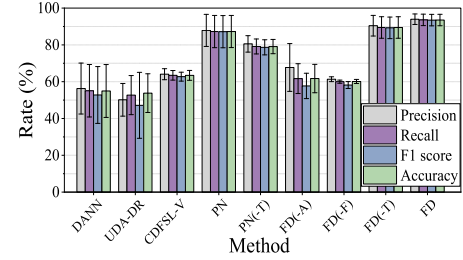


Fig. 18. The performance of cross-subject for different methods.

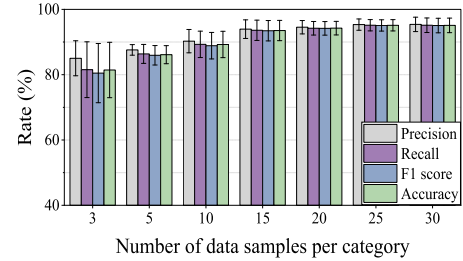


Fig. 19. The performance of different numbers of data samples per category.

Additionally, Fig. 18 shows the cross-subject depressive episode monitoring results using the unsupervised domain adaptation methods DANN [42], UDA-DR, the few-shot learning methods CDFSL-V [46], PN [47], FSL-DEM, as well as ablation versions of PN and FSL-DEM. The results show that the performance of unsupervised domain adaptation has largely improved compared to using only LCMTN for cross-subject performance. The FSL-DEM exhibits the best performance compared to the other two few-shot learning methods. Furthermore, we conduct an additional ablation study on FSL-DEM. From the comparison between PN and PN(-T), as well as FD and FD(-T), utilizing a subject-based task generation strategy can enhance the cross-subject performance of FSL-DEM. Comparing FD and FD(-F), it is evident that the fine-tuning module is indispensable. The pre-training approach based on meta-learning is superior to the conventional training method. Comparing FD and FD(-A), solely utilizing readily available non-episode samples does not improve performance.

2) *Analysis of Performance Influencing Factors:* To better select appropriate hyperparameters, we also evaluate four influencing factors affecting the cross-subject performance of FSL-DEM: the number of internal training, the number of data samples per category, the number of fine-tuning iterations, and the number of source domain subjects.

The number of data samples per category: In the fine-tuning phase, we provide 3, 5, 10, 15, 20, 25, and 30 labeled samples for each category. For fine-tuning, the number of fine-tuning iterations is fixed. As shown in Fig. 19, the performance metrics of precision, recall, F1 score, and accuracy increase as the number of samples provided by the target subject increases. When providing 3 samples for each category, the accuracy can reach 81.43%. After providing 15 samples for each category, the performance metrics tend to stabilize. During the experimental phase of collecting data on depressive episode, it was found that each episode of patients lasted at least 15 seconds. Therefore,

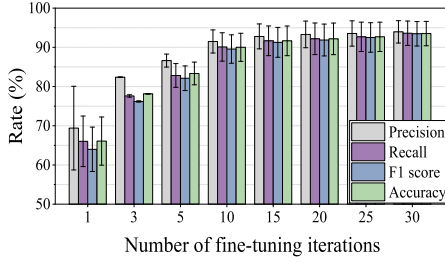


Fig. 20. The performance of different numbers of fine-tuning iterations.

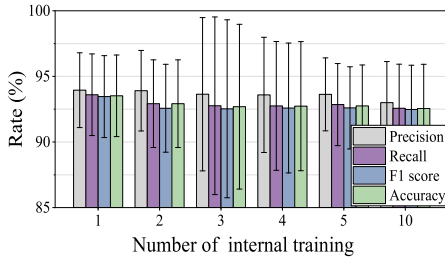


Fig. 21. The performance of different numbers of internal training.

considering both the cost of target subject data and the performance of depressive episode monitoring, setting 15 labeled samples for each category as the minimum required data from the target subject is advisable.

The number of fine-tuning iterations: In the fine-tuning phase, we fix the number of samples provided by the target subject for each category at 15 and set fine-tuning iterations to 1, 3, 5, 10, 15, 20, 25, and 30, respectively. Observing from Fig. 20, it can be noted that the performance metrics of precision, recall, F1 score, and accuracy exhibit a rapid increase followed by a gradual stabilization as the number of fine-tuning iterations increases. After 15 fine-tuning iterations, the performance improvement is small. Specifically, the accuracy rates are 91.66%, 92.17%, 92.68%, and 93.52% for fine-tuning iterations of 15, 20, 25, and 30, respectively. Therefore, we fixed the number of fine-tunings to 20 times.

The number of internal training: In the meta-learning phase, the number of internal training not only affects the training duration but also influences the model performance. Specifically, we conduct performance evaluations for depressive episode monitoring by setting the internal training iterations to 1, 2, 3, 4, 5, and 10. As shown in Fig. 21, precision, recall, F1 score, and accuracy all decrease as the number of internal training iterations increases, but the decrease is very small. This is because the increase in the number of internal training iterations leads to overfitting of the model on specific tasks. Among them, when the number of internal training iterations is 1, all performance metrics are optimal. Therefore, choosing an internal training iteration of 1 not only reduces meta-learning training time but also results in better performance in depressive episode monitoring.

The number of source domain subjects: The diversity of subjects is closely related to the generalization performance of the model. In the meta-learning phase, apart from each individual subject, we randomly select 3, 5, 7, 10, 15, 17 subjects from the remaining subjects as source domain subjects. As shown in

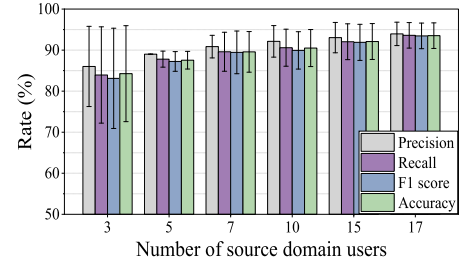


Fig. 22. The performance of different numbers of source domain subjects.

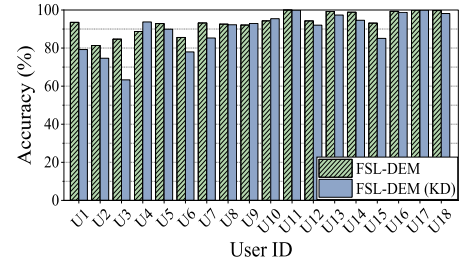


Fig. 23. The results of FSL-DEM and FSL-DEM (KD).

Fig. 22, as the number of source domain subjects increases, the performance of depressive episode monitoring gradually improves. With only three source domain subjects, the accuracy of depressive episode monitoring reaches 84.26%, which is a 29.8% improvement compared to the cross-subject accuracy of the LCMTN. From having 3 source domain subjects to 15, the accuracy increases from 84.26% to 92.1%. Further increasing to include all source domain subjects results in only a 1.42% improvement in accuracy. Even if the number of source domain subjects is further increased, the room for performance improvement is extremely limited while the data cost continues to increase.

3) Impact of Knowledge Distillation: A knowledge distillation module is designed to reduce computational complexity and memory usage, while ensuring the performance of depressive episode monitoring. The distillation degree τ refers to a parameter that quantifies the intensity of knowledge distillation during the knowledge distillation process. In this regard, we use FSL-DEM and FSL-DEM (KD) with a distillation degree of 10 for the encoder. As shown in Fig. 23, the accuracy after distillation for most subjects is slightly lower than the accuracy before distillation. The average accuracy before and after distillation is 93.52% and 90.53, with a difference of only 2.99%. Additionally, subjects U4, U9, and U10 have improved accuracy through knowledge distillation, indicating that knowledge distillation can optimize performance for certain subjects. Furthermore, we evaluate the impact of knowledge distillation at different degrees on the performance of depressive episode monitoring. We respectively set the degrees of knowledge distillation for the encoder to 5, 10, 20, and 40, and obtain the results shown in Fig. 24. It can be observed that as the degree of knowledge distillation increases, various performance metrics show a gradual decrease. It is worth noting that although the degree of knowledge distillation increases from 5 to 40, the performance has declined slightly, but not significantly. Therefore, when

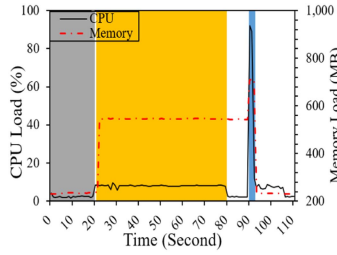


Fig. 28. CPU and memory load.

TABLE V
THE COMPARISON OF OUR SELF-DESIGNED HARDWARE
WITH COMMERCIAL DEVICES

Device	Size (cm)	Mass (g)	Price (\$)	Data Interface
Silme W20 ²	2.05 × 6.5 × 1.25	27.5	538.56	not available
Fitbit Versa 4 ³	4.05 × 4.05 × 1.12	37.64	120	not available
Empatica E4 ⁴	4.4 × 4 × 19	25	1080	pay 1080 \$ for 3-year data access
Ours (Ver 1.0)	7 × 4 × 2.5	120	23.56	available
Ours (Ver 2.0)	3.8 × 3.8 × 1.1	28.6	18.57	available

respectively. These results highlight the student model's suitability for edge deployment, emphasizing the benefits of knowledge distillation for resource-constrained devices.

3) *CPU and Memory Load*: In this experiment, we demonstrate the dynamics of the CPU and memory footprint of a student model deployed on an edge device for acquisition. All other applications and services are closed to ensure accurate measurements of CPU and memory load. Android Profiler in USB debug mode is used to monitor real-time CPU and memory usage for a duration of 110 seconds. As shown in Fig. 28, the results are mainly divided into three phases. In the initial stage of the silent state (highlighted in gray), CPU usage remains minimal, and memory usage stabilizes at around 240 MB. During the data acquisition and visualization phase (highlighted in yellow), memory usage rises to approximately 550 MB, while CPU load increases slightly to 10%. We acquire 60 seconds of physiological signaling data during this phase. In the data processing and model prediction phase (highlighted in blue), which lasts about 3 seconds, resource consumption surges, with memory usage reaching about 720 MB and CPU utilization peaking at 92%, highlighting the computational intensity of the predictive modeling process.

VI. DISCUSSION AND FUTURE WORK

In this section, we discuss some important details of our device and the future work.

A. Wearable Devices

As shown in Table V, we compare our self-designed hardware with existing commercial devices in terms of size, weight, cost, and data interface. We can see that these commercial devices are much more expensive than ours or do not provide access to sensor data. This is why we need to design such a device by ourselves instead of using commercial ones. Specifically, the cost of our device (Ver 1.0) encompasses not only the four sensors detailed in Table I but also expenditures of the PCB circuit board (0.70 \$) and 3D-printed casing (6.382 \$). The

entire hardware is encapsulated within a custom-designed 3D printed enclosure with dimensions of 7 cm × 4 cm × 2.5 cm. We can also see that our device (Ver 1.0) is larger and heavier than other commercial devices, as our design incorporates several physiological sensors through an embedded system approach. However, by enhancing the integration density of hardware circuits and optimizing the lightweight design of the enclosure, we have enabled the device (Ver 2.0) to achieve more compact size and reduced weight, as shown in Table V.

It is worth noting that, despite our device's affordability, it still achieves accurate depression detection recognition for two key reasons. On the one hand, as detailed in Table I, our device supports much higher sampling rates for four kinds of sensors than commercial devices such as Empatica E4 with a maximum sampling rate of 64 Hz. The high sampling rates allow us to capture higher resolution and more detailed temporal signals. On the other hand, we have proposed novel methods to address the issue of human heterogeneity and enhance cross-subject recognition performance.

B. Real-World Adaptation

Although disease recognition and monitoring increasingly depend on data from real-world, dynamic environments in the current era of ubiquitous computing, our study remains limited in adapting to this trend.

(1) *Environmental inconsistencies and noise interference*: The differences of data-collection environments create significant challenges. Real-world deployment faces unpredictable factors, including environmental noise, device variability, and patients' daily routines. At present, we require subjects to remain in a resting state during data collection, while noises generated by human motion artifacts inevitably affect the data quality of the physiological signals. In the future, we will explore reasonable artifact noise processing algorithms (e.g., wavelet algorithm) to mitigate motion artifacts and enhance robustness under naturalistic conditions.

(2) *Model generalizability and adaptability*: Model adaptability is also an important issue. Models optimized for specific tasks and data distributions in controlled environments may struggle with the diverse and dynamic conditions and requirements in real-world settings. For example, a model trained to extract depression-related features from particular physiological signals may not promptly recognize and adjust to the wider range of symptoms and behaviors exhibited by patients in the wild, potentially resulting in suboptimal effectiveness. We will continue to improve the *DepGuard* by enhancing its flexibility and adaptability to a broader population and more dynamic conditions, potentially drawing on recent advances in multi-modal federated learning system [53].

(3) *Longitudinal validation*: Due to the short hospitalization time of depressed patients, our work lacks a cross-session test to validate longitudinal performance across time. This limitation hinders our understanding of model consistency and reliability over extended periods. In future work, we aim to distribute the hardware and software to experimenters, enabling them to collect richer data and achieve real-time monitoring in their daily lives. This will facilitate long-term evaluation and validation.

VII. CONCLUSION

In summary, we present *DepGuard*, a wearable depression assessment system that integrates a wrist-worn device with a mobile app for both depression recognition and real-time monitoring of depressive episodes. The system is designed for potential clinical application to help alleviate the workload of healthcare professionals. To address cross-domain generalization challenges, we employ unsupervised domain adaptation and few-shot learning techniques. Extensive experiments validate the effectiveness of our approach. Results indicate a strong correlation between heart rate and depression recognition, as well as between skin temperature and depressive episode detection, further reinforcing findings in psychiatric research.

REFERENCES

- [1] W. H. Organization et al., *Depressive Disorder (Depression) World Health Organization*. Geneva, Switzerland: WHO, 2023.
- [2] D. American Psychiatric Association et al., *Diagnostic and Statistical Manual of Mental Disorders: DSM-5*, vol. 5. Washington, DC, USA: American Psychiatric Association, 2013.
- [3] H. Herrman et al., "Time for united action on depression: A lancet-world psychiatric association commission," *Lancet*, vol. 399, no. 10328, pp. 957–1022, 2022.
- [4] X. Xu et al., "Mental-LLM: Leveraging large language models for mental health prediction via online text data," in *Proc. ACM Interactive Mobile Wearable Ubiquitous Technol.*, 2024, pp. 1–32.
- [5] X. Xu et al., "Leveraging collaborative-filtering for personalized behavior modeling: A case study of depression detection among college students," in *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 5, no. 1, pp. 1–27, 2021.
- [6] Y. Rykov et al., "Digital biomarkers for depression screening with wearable devices: Cross-sectional study with machine learning modeling," *JMIR mHealth uHealth*, vol. 9, no. 10, 2021, Art. no. e24872.
- [7] M. Tlachac, R. Flores, M. Reisch, K. Houskeeper, and E. A. Rundensteiner, "DepreST-CAT: Retrospective smartphone call and text logs collected during the COVID-19 pandemic to screen for mental illnesses," in *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 6, no. 2, pp. 1–32, 2022.
- [8] M. A. Wani, M. A. ELAffendi, K. A. Shakil, A. S. Imran, and A. A. Abd El-Latif, "Depression screening in humans with AI and deep learning techniques," *IEEE Trans. Computat. Social Syst.*, vol. 10, no. 4, pp. 2074–2089, Aug. 2023.
- [9] C. Lin et al., "Sensemood: Depression detection on social media," in *Proc. 2020 Int. Conf. Multimedia Retrieval*, 2020, pp. 407–411.
- [10] S. Nepal et al., "Capturing the college experience: A four-year mobile sensing study of mental health, resilience and behavior of college students during the pandemic," in *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 8, no. 1, pp. 1–37, 2024.
- [11] X. Xu et al., "GLOBEM: Cross-dataset generalization of longitudinal human behavior modeling," in *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 6, no. 4, pp. 1–34, 2023.
- [12] K. Opoku Asare, Y. Terhorst, J. Vega, E. Peltonen, E. Lagerspetz, and D. Ferreira, "Predicting depression from smartphone behavioral markers using machine learning methods, hyperparameter optimization, and feature importance analysis: Exploratory study," *JMIR mHealth uHealth*, vol. 9, no. 7, 2021, Art. no. e26540.
- [13] Y. Pan et al., "Spatial-temporal attention network for depression recognition from facial videos," *Expert Syst. Appl.*, vol. 237, 2024, Art. no. 121410.
- [14] C. Koch, M. Wilhelm, S. Salzmann, W. Rief, and F. Euteneuer, "A meta-analysis of heart rate variability in major depression," *Psychol. Med.*, vol. 49, no. 12, pp. 1948–1957, 2019.
- [15] S. Saganowski, B. Perz, A. G. Polak, and P. Kazienko, "Emotion recognition for everyday life using physiological signals from wearables: A systematic literature review," *IEEE Trans. Affective Comput.*, vol. 14, no. 3, pp. 1876–1897, Third Quarter, 2023.
- [16] G. Tasci et al., "Automated accurate detection of depression using twin Pascal's triangles lattice pattern with EEG signals," *Knowl.-Based Syst.*, vol. 260, 2023, Art. no. 110190.
- [17] H. Lin et al., "MDD-TSVM: A novel semisupervised-based method for major depressive disorder detection using electroencephalogram signals," *Comput. Biol. Med.*, vol. 140, 2022, Art. no. 105039.
- [18] S. Byun et al., "Detection of major depressive disorder from linear and nonlinear heart rate variability features during mental task protocol," *Comput. Biol. Med.*, vol. 112, 2019, Art. no. 103381.
- [19] D. Wang, J. Weng, Y. Zou, and K. Wu, "Emotracer: A wearable physiological and psychological monitoring system with multi-modal sensors," in *Proc. 2022 ACM Int. Joint Conf. Pervasive Ubiquitous Comput.-2022 ACM Int. Symp. Wearable Comput.*, 2022, pp. 444–449.
- [20] Y. Zou et al., "Research on wearable emotion recognition based on multi-source domain adversarial transfer learning (in chinese)," *Chin. J. Comput.*, vol. 47, no. 2, pp. 266–286, 2024.
- [21] S. Hassantabar, J. Zhang, H. Yin, and N. K. Jha, "MHDeep: Mental health disorder detection system based on wearable sensors and artificial neural networks," *ACM Trans. Embedded Comput. Syst.*, vol. 21, no. 6, pp. 1–22, 2022.
- [22] Y. Jiao et al., "Feasibility study for detection of mental stress and depression using pulse rate variability metrics via various durations," *Biomed. Signal Process. Control*, vol. 79, 2023, Art. no. 104145.
- [23] Y. Hu, J. Chen, J. Chen, W. Wang, S. Zhao, and X. Hu, "An ensemble classification model for depression based on wearable device sleep data," *IEEE J. Biomed. Health Inform.*, vol. 28, no. 5, pp. 2602–2612, May 2024.
- [24] W. Kuang, S. Jin, D. Wang, Y. Zheng, Y. Zou, and K. Wu, "EmoTracer: A user-independent wearable emotion tracer with multi-source physiological sensors based on few-shot learning," in *Proc. 2024 IEEE Int. Conf. Bioinf. Biomed.*, 2024, pp. 2680–2687.
- [25] A. Ahmed, J. Ramesh, S. Ganguly, R. Aburukba, A. Sagahyroon, and F. Aloul, "Evaluating multimodal wearable sensors for quantifying affective states and depression with neural networks," *IEEE Sensors J.*, vol. 23, no. 19, pp. 22788–22802, Oct. 2023.
- [26] G. Sharma, A. Parashar, and A. M. Joshi, "DepHNN: A novel hybrid neural network for electroencephalogram (EEG)-based screening of depression," *Biomed. Signal Process. Control*, vol. 66, 2021, Art. no. 102393.
- [27] F. Zhang, M. Wang, J. Qin, Y. Zhao, X. Sun, and W. Wen, "Depression recognition based on electrocardiogram," in *Proc. IEEE Int. Conf. Comput. Commun. Syst.*, 2023, pp. 1–5.
- [28] C.-C. Huang et al., "MEG-based classification and grad-cam visualization for major depressive and bipolar disorders with semi-CNN," in *Proc. 44th Annu. Int. Conf. IEEE Eng. Med. Biol. Soc.*, 2022, pp. 1823–1826.
- [29] K. Shao, Y. Hao, L. Hu, X. Zong, and M. Chen, "Data augmentation and pseudo-sequence of fNIRS for depression recognition," in *Proc. 2023 IEEE Int. Conf. Bioinf. Biomed.*, 2023, pp. 2223–2226.
- [30] M. Mousavian, J. Chen, Z. Traylor, and S. Greening, "Depression detection from sMRI and rs-fMRI images using machine learning," *J. Intell. Inf. Syst.*, vol. 57, pp. 395–418, 2021.
- [31] S. Pange and V. Pawar, "Depression analysis based on EEG and ECG signals," in *Proc. IEEE Int. Conf. Emerg. Technol.*, 2023, pp. 1–6.
- [32] Y. Tazawa et al., "Evaluating depression with multimodal wristband-type wearable device: Screening and assessing patient severity utilizing machine-learning," *Heliyon*, vol. 6, no. 2, 2020, Art. no. 20219964001.
- [33] I. Moshe et al., "Predicting symptoms of depression and anxiety using smartphone and wearable data," *Front. Psychiatry*, vol. 12, 2021, Art. no. 625247.
- [34] Y. Fang, J. Wu, Q. Wang, S. Qiu, A. Bozoki, and M. Liu, "Source-free collaborative domain adaptation via multi-perspective feature enrichment for functional MRI analysis," *Pattern Recognit.*, vol. 157, 2025, Art. no. 110912.
- [35] T. Chen, Y. Guo, S. Hao, and R. Hong, "Semi-supervised domain adaptation for major depressive disorder detection," *IEEE Trans. Multimedia*, vol. 26, pp. 3567–3579, 2024.
- [36] Z. Zhang, Q. Meng, L. Jin, H. Wang, and H. Hou, "A novel EEG-based graph convolution network for depression detection: Incorporating secondary subject partitioning and attention mechanism," *Expert Syst. Appl.*, vol. 239, 2024, Art. no. 122356.
- [37] Y. Fang, M. Wang, G. G. Potter, and M. Liu, "Unsupervised cross-domain functional mri adaptation for automated major depressive disorder identification," *Med. Image Anal.*, vol. 84, 2023, Art. no. 102707.
- [38] Z. Zhang, C. Xu, L. Jin, H. Hou, and Q. Meng, "A depression level classification model based on EEG and GCN network with domain generalization," in *Proc. IEEE 43rd Chin. Control Conf.*, 2024, pp. 8582–8587.
- [39] A. Vaswani et al., "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 6000–6010.
- [40] C. Grady, "Institutional review boards: Purpose and challenges," *Chest*, vol. 148, no. 5, pp. 1148–1155, 2015.

- [41] D. P. Kingma, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*.
- [42] Y. Ganin et al., "Domain-adversarial training of neural networks," *J. Mach. Learn. Res.*, vol. 17, no. 59, pp. 1–35, 2016.
- [43] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 7167–7176.
- [44] J. Wang, J. Chen, J. Lin, L. Sigal, and C. W. de Silva, "Discriminative feature alignment: Improving transferability of unsupervised domain adaptation by gaussian-guided latent alignment," *Pattern Recognit.*, vol. 116, 2021, Art. no. 107943.
- [45] R. Hu, L. Chen, S. Miao, and X. Tang, "SWL-Adapt: An unsupervised domain adaptation model with sample weight learning for cross-user wearable human activity recognition," in *Proc. AAAI Conf. Artif. Intell.*, 2023, pp. 6012–6020.
- [46] S. Samarasinghe, M. N. Rizve, N. Kardan, and M. Shah, "CDFSL-V: Cross-domain few-shot learning for videos," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2023, pp. 11 643–11 652.
- [47] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 4080–4090.
- [48] M. Kobayashi, G. Sun, T. Shinba, T. Matsui, and T. Kirimoto, "Development of a mental disorder screening system using support vector machine for classification of heart rate variability measured from single-lead electrocardiography," in *Proc. IEEE Sensors Appl. Symp.*, 2019, pp. 1–6.
- [49] L. Yang, Y. Wang, X. Zhu, X. Yang, and C. Zheng, "A gated temporal-separable attention network for EEG-based depression recognition," *Comput. Biol. Med.*, vol. 157, 2023, Art. no. 106782.
- [50] K. Shao, R. Wang, Y. Hao, L. Hu, M. Chen, and H. Arno Jacobsen, "Multimodal physiological signals representation learning via multiscale contrasting for depression recognition," in *Proc. 32nd ACM Int. Conf. Multimedia*, 2024, pp. 5692–5701.
- [51] C. Fu et al., "M 3 ADD: A novel benchmark for physiology signal-based automatic depression detection with multimodal multitask multievent framework," in *Proc. 2025 IEEE Int. Conf. Acoust. Speech Signal Process.*, 2025, pp. 1–5.
- [52] X. Shui, H. Xu, S. Tan, and D. Zhang, "Depression recognition using daily wearable-derived physiological data," *Sensors*, vol. 25, no. 2, 2025, Art. no. 567.
- [53] X. Ouyang et al., "ADMarker: A multi-modal federated learning system for monitoring digital biomarkers of alzheimer's disease," in *Proc. 30th Annu. Int. Conf. Mobile Comput. Netw.*, 2024, pp. 404–419.



Yufei Zhang received the BS and MS degrees from the School of Computer and Information Engineering, Shanghai Polytechnic University (SSPU), in 2021 and 2024, respectively. He is currently working toward the PhD degree with the College of Computer Science and Software Engineering, Shenzhen University. His research interests include affective computing, mobile computing, and spatial-temporal data analysis.



Shuo Jin received the master's degree from the College of Computer Science and Software Engineering, Shenzhen University in 2025. He is currently an algorithm engineer with GYENNO Technologies. He participated in this work during working toward the master's degree with Shenzhen University. His research interests focus on mobile computing and affective computing.



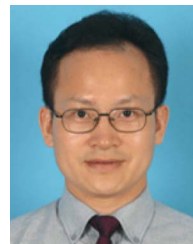
Wenting Kuang is currently working toward the master's degree with the College of Computer Science and Software Engineering, Shenzhen University, Shenzhen, China. Her research interests covers affective computing and intelligent sensing.



Yuda Zheng is currently working toward the bachelor's degree with the College of Computer Science and Software Engineering, Shenzhen University, China. His research focuses on mobile application development and affective computing.



Qifeng Song received the MS degree from the College of Computer Science and Software Engineering, Shenzhen University, in 2024. He is currently an algorithm engineer with LEPU MEDICAL. He participated in this work during working toward the master's degree with Shenzhen University. His research interests include affective computing and mobile computing.



Changhe Fan is currently a chief physician from the Department of Psychiatry, Guangdong Second Provincial General Hospital. His research interests include biological psychiatry and behavioral medicine.



Yongpan Zou (Member, IEEE) received the PhD degree from the Department of Computer Science and Engineering (CSE), Hong Kong University of Science and Technology, in 2017. He is currently an associate professor with the College of Computer Science and Software Engineering, Shenzhen University. His research interests include ubiquitous sensing, mobile computing, and human-computer interaction.



Victor C. M. Leung (Life Fellow, IEEE) received the PhD degree in electrical engineering from the University of British Columbia, in 1981. He currently is the dean with the Artificial Intelligence Research Institute, Shenzhen MSU-BIT University. His research focuses on wireless networks and mobile systems with more than 60,000 citations. He has received numerous accolades, such as the APEBC Gold Medal, NSERC Postgraduate Scholarships, and IEEE Vancouver Section Centennial Award.



Kaishun Wu (Fellow, IEEE) received the PhD degree from the Hong Kong University of Science and Technology, Hong Kong, in 2011. He is currently a professor in information hub with the Hong Kong University of Science and Technology (Guangzhou). His research interests include wireless communications and mobile computing. He won several best paper awards of international conferences, such as IEEE Globecom 2012 and IEEE MASS 2014.