# MHDeep: Mental Health Disorder Detection System based on Body-Area and Deep Neural Networks

Shayan Hassantabar, Joe Zhang, Hongxu Yin, and Niraj K. Jha, *Fellow, IEEE*

**Abstract**—Mental health problems impact quality of life of millions of people around the world. However, diagnosis of mental health disorders is a challenging problem that often relies on self-reporting by patients about their behavioral patterns and social interactions. Therefore, there is a need for new strategies for diagnosis and daily monitoring of mental health conditions. The recent introduction of body-area networks consisting of a plethora of accurate sensors embedded in smartwatches and smartphones and edge-compatible deep neural networks (DNNs) points towards a possible solution. Such wearable medical sensors (WMSs) enable continuous monitoring of physiological signals in a passive and non-invasive manner. However, disease diagnosis based on WMSs and DNNs, and their deployment on edge devices, such as smartphones, remains a challenging problem. To this end, we propose a framework called MHDeep that utilizes commercially available WMSs and efficient DNN models to diagnose three important mental health disorders: schizoaffective, major depressive, and bipolar. MHDeep uses eight different categories of data obtained from sensors integrated in a smartwatch and smartphone. These categories include various physiological signals and additional information on motion patterns and environmental variables related to the wearer. MHDeep eliminates the need for manual feature engineering by directly operating on the data streams obtained from participants. Since the amount of the data is limited, MHDeep uses a synthetic data generation module to augment real data with synthetic data drawn from the same probability distribution. We use the synthetic dataset to pre-train the weights of the DNN models, thus imposing a prior on the weights. We use a grow-and-prune DNN synthesis approach to learn both the architecture and weights during the training process. We use three different data partitions to evaluate the MHDeep models trained with data collected from 74 individuals. We conduct two types of evaluations: at the data instance level and at the patient level. MHDeep achieves an average test accuracy, across the three data partitions, of 90.4%, 87.3%, and 82.4%, respectively, for classifications between healthy and schizoaffective disorder instances, healthy and major depressive disorder instances, and healthy and bipolar disorder instances. At the patient level, MHDeep DNNs achieve an accuracy of 100%, 100%, and 90.0% for the three mental health disorders, respectively, based on inference that uses 40, 16, and 22 minutes of data from each patient.

**Index Terms**—Body-area network; deep neural network; disease diagnosis; health monitoring; mental health disorders; wearable medical sensors; synthetic data generation.

<div align="center">✦</div>

## 1 INTRODUCTION

Mental health problems impact around 20% of the world population [1]. They may negatively affect a person's mind, emotions, behavior, and even physical health. Mental health issues may include various disorders like bipolar, depression, schizophrenia, and attention-deficit hyperactivity, to name but a few. These disorders do not only affect adults, but children and adolescents may suffer from them as well [2]. Moreover, patients with serious mental health issues have a higher risk of morbidity due to physical health problems. Depression, for example, can lead to disability and increases the risk of suicidal thoughts and attempts.

In order to understand the mental health condition of the patient and provide suitable patient care, early detection is essential. However, this remains a public health challenge. While many other diseases can be diagnosed based on specific medical tests and laboratory measurements, detection of mental health problems mainly relies on self-reports and responses to specific questionnaires designed for identifying certain patterns of behavior and social interactions. Hence, to address this challenge, novel detection strategies are needed.

There has been recent interest in employing machine learning (ML) to detect mental health conditions [3]. Neural networks are popular machine learning models that use nonlinear computations to make inferences from large datasets. Thus, they have started being deployed in the smart healthcare domain [4]–[9].

In previous studies, two main data sources for deep learning based analysis of mental health have been clinical data and social media usage data. The former includes studies that use neuro-image data for detecting various mental health disorders [10], electroencephalogram (EEG) data to study brain disorders [11], and analysis of electronic health records (EHR) to study mental health problems [12]. Moreover, social media usage patterns have been used to predict personal traits of the user. As a result, several recent works focus on exploiting such patterns to detect psychiatric illness [13].

Although the above works have demonstrated the promise of using machine learning in identifying mental health disorders, daily mental health monitoring is still a challenge. Since mental health condition treatment delays may lead to negative outcomes, potentially even loss of life, it is desirable to have immediate and pervasive mental health detection. This is the motivation behind our mental health de-

tection system, MHDeep. As shown in Fig. 1, MHDeep relies on physiological data collected using various WMSs. WMSs can be used to continuously monitor the physiological signals of the wearer throughout the day. This enables constant tracking of the health conditions of the user. MHDeep uses various sensors embedded in smartwatches and smartphones. For training purposes, the collected physiological data are processed to obtain a comprehensive dataset. MHDeep combines data from WMSs with the inference capabilities of deep neural networks (DNNs) to directly extract mental health condition from the physiological signals. These inferences can be communicated to a health server that is accessible to the physician. This has the potential to enhance to ability of the physician to intervene quickly when mental health conditions deteriorate.

Difficulty of data collection and labeling limits the amount of available data. Hence, MHDeep uses synthetic data drawn from the same probability distribution as real data to augment the dataset. It also leverages a grow-and-prune DNN synthesis approach [14], [15] to train accurate and computationally efficient neural network models to detect the mental health condition of the user.

The major contributions of this article are summarized next.

- We demonstrate an easy-to-use, accurate, and pervasive mental health disorder detection system, called MHDeep. MHDeep combines physiological signals collected from WMSs with the prediction power of DNNs to detect three main mental health disorders: major depressive, schizoaffective, and bipolar. Unlike many other approaches for detecting mental health problems, MHDeep does not rely on any self-reports from the user.
- We do an extensive search to extract the most appropriate set of data categories for detecting each of the three mental health disorders.
- MHdeep relies on a synthetic data generation module to alleviate the concerns arising from unavailability of large datasets. It uses a grow-and-prune DNN synthesis approach to improve the accuracy of the DNNs while reducing their computational costs.
- We demonstrate the performance, accuracy, and feasibility of MHDeep through extensive evaluations.

The rest of the article is organized as follows. Section 2 presents background information on various works related to MHDeep. Section 3 explains the MHDeep framework thoroughly. Section 4 provides implementation details. Section 5 presents experimental evaluations. Section 6 provides a short discussion related to this work. Finally, Section 7 concludes the article.

## 2 BACKGROUND

In this section, we first provide background information on various mental health disorders and how they affect patient lives. Next, we discuss various methods for identifying mental health conditions based on machine learning. We also discuss some of the related work on synthesizing efficient neural network models. Finally, we discuss WMSs and their applications to various disease diagnosis frameworks.
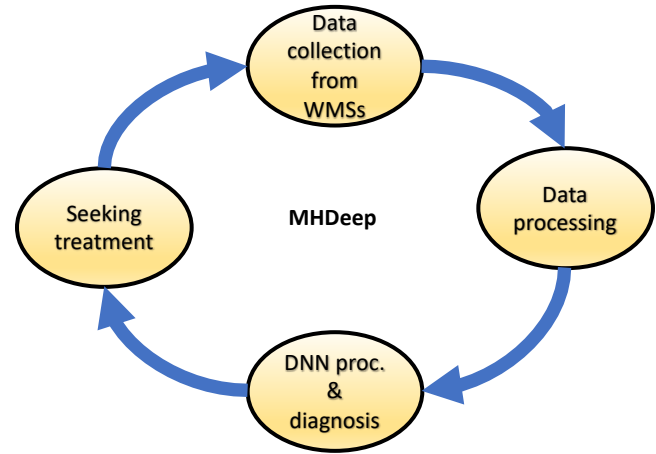


Fig. 1. MHDeep mental health disorder detection system

### 2.1 Mental health disorders and their impact

Mental health conditions can affect a patient's thinking, feeling, and behavior. They may have a deep impact on the daily life of the person and affect their ability to adequately perform in society. There are hundreds of different mental illnesses [16]. We discuss the three mental health disorders that we target in this work: bipolar, major depressive, and schizoaffective.

Bipolar disorder can cause a dramatic shift in a person's mood, energy, and behavior. It is characterized by experiences of alternating episodes of manic and depressed states. Major depressive disorder may present different symptoms like loss of interest, sleep disturbance, change in appetite, and feeling of fatigue. Schizoaffective disorder is characterized by various symptoms of schizophrenia such as episodes of hallucinations and delusions. It may also present other symptoms such as disorganized thinking, depressed mood, and manic behavior.

Apart from various conditions that these mental health disorders may cause, stereotypes related to mental health seem to still be widely prevalent in society, not just among uninformed people but even among well-trained professionals [17]. These stereotypes often lead to social and employment discrimination [18] and poor treatment of physical health problems [19].

### 2.2 Deep learning for mental health

Deep learning has been recently used to better understand and detect mental health problems. Deep learning approaches have been applied to various types of data: mainly clinical data and social media usage data [20]. The three types of clinical data used in these works are neuroimage data, EEG data, and EHR data.

Several studies demonstrate the effectiveness of neuroimages in detecting neuropsychiatric disorders [10]. Two types of neuroimage data used in such works are functional magnetic resonance imaging (fMRI) and structural magnetic resonance imaging (sMRI). fMRI measures brain activity by monitoring blood oxygenation and flow in response to neural activity. sMRI examines the anatomy and pathology of the brain. Deep belief networks have been used to detect the

presence of attention-deficit hyperactivity disorder (ADHD) using fMRI and sMRI data [21], [22]. These data types have also been used to detect schizophrenia [23], [24]. Depression has been detected using time-series fMRI data using convolutional neural networks (CNNs) and autoencoders [25]. EEG is another source of data for studying brain disorders. For example, CNN based feature extraction from EEG data has been used to detect depression [11]. EHR is a collection of patient-centered records and includes both structural data such as laboratory reports and unstructured data such as clinical and discharge notes. Since EHR data is a collection of longitudinal records, recurrent neural networks (RNNs) have been used to distill information from them. Pham et al. [26] use RNN architectures to predict future outcomes of depressive episodes. Unstructured clinical notes have been analyzed with deep learning-based models to detect depression [12].

Social media usage data have also proved their usefulness in identifying psychiatric illness. Birnbaum et al. [13] investigate Facebook messages and patterns of sharing images on social media to distinguish among healthy individuals, individuals with a schizophrenia spectrum disorder, and individuals with mood disorders. Other works have used DNNs with textual data and image data shared on social media platforms to detect stress [27], depression [28], and risk of suicide [29].

### 2.3 Efficient neural network synthesis

Next, we summarize the main approaches to synthesis of compact DNN models. Conventional synthesis methods are based on the use of efficient building blocks. For example, MobileNetV2 [30] leverages inverted residual blocks to reduce model size and computations significantly. Wu et al. [31] use shift-based operations rather than convolution layers to significantly reduce computational costs of the model. The main drawback of such approaches is the need for considerable design insight and trial-and-error process to design such efficient building blocks. Network compression is another approach to design of efficient models. It removes the need for design insights. Network pruning is a widely used method that eliminates weights or filters that do not enhance model performance. Han et al. [32] have shown the effectiveness of pruning in removing redundancy in CNNs and multilayer-perceptron architectures. Grow-and-prune DNN synthesis uses network growth followed by network pruning in an iterative process to improve model performance while ensuring its compactness [14], [15].

Another recent approach relies on the use of reinforcement learning (RL) to search for DNN architectures in an automated flow. It is known as neural architecture search (NAS) [33]. NAS generally uses a controller, e.g., an RNN, to iteratively generate candidate architectures in the search process. The RL controller is improved based on candidate performance. As an example, MnasNet [34] uses an RL-based approach to develop efficient DNNs for mobile platforms. However, the downside of the RL-based NAS approach is that it is computationally intensive. FBNet [35] uses the Gumbel softmax function to optimize weights and connections using a single objective function. NEAT [36] uses evolutionary algorithms to generate optimized and increasingly complex architectures over multiple generations. Combining efficient evolutionary search algorithms with various performance predictors, e.g., for accuracy, energy, and latency, is another approach for synthesizing accurate yet compact CNNs/DNNs [37], [38].

### 2.4 Wearable medical sensors

Due to recent developments in low-power sensor design and efficient wireless communications, battery-powered WMSs are becoming ubiquitous. More than 123 million WMSs were sold worldwide in 2018. This number is projected to grow to 1 billion by the end of 2022 [39]. WMSs can track different aspects of human health including heart rate, body/skin temperature, respiration rate, blood pressure, EEG, electrocardiogram (ECG), and Galvanic skin response (GSR) [40]. Furthermore, the number of physiological signals that can be measured using WMSs keeps growing every year.

WMSs have begun to be used in many smart healthcare applications. CodeBlue [41] is a sensor network that collects vital health signs and transmits them to the healthcare provider. MobiHealth [42] is a WMS based body-area network (BAN) that realizes an end-to-end mobile health monitoring platform. Yin et al. [7] use WMSs for pervasive diagnosis of Type-I and Type-II diabetes. CovidDeep [4] is a WMS-based framework for quick detection of SARS-CoV-2/COVID-19.

For data collection in MHDeep, we use an Empatica E4 smartwatch [43] to record a subset of patient's physiological signals. It is a wearable wireless device designed for comfortable, continuous, and real-time data acquisition. We also use a smartphone to simultaneously record signals related to motion information and environmental variables. Since the DNNs developed for diagnosing various mental health conditions can reside on the smartphone, use of a smartwatch/smartphone based BAN can enable accurate, yet convenient, disease diagnosis and continuous healthcare monitoring.

## 3 METHODOLOGY

In this section we describe various parts of the MHDeep framework. First, we give an overview of our approach. Then, we discuss the data collection and preparation process, synthetic data generation, and grow-and-prune DNN synthesis.

### 3.1 The MHDeep framework

We illustrate the MHDeep framework in Fig. 2. The input data are derived from the physiological signals collected using various WMSs in the smartwatch and smartphone in a non-invasive, passive, and efficient manner. The list of collected data streams include GSR, skin temperature (ST), inter-beat interval (IBI), and 3-way acceleration from the smartwatch. In addition, some information related to motion patterns of the user and ambient information are collected using smartphone sensors. This includes ambient temperature, gravity, acceleration, and angular velocity. After sensor data collection, the collected signals are synchronized, aggregated, and merged into a comprehensive data input for subsequent analysis. To enhance the accuracy of subsequent analysis and

improve noise tolerance, we normalize the data. The process of data collection and preparation is discussed in more detail in Section 3.2. When the size of the training dataset is small, it can be useful to generate a synthetic dataset from the same probability distribution as the real training dataset. MHDeep leverages Gaussian mixture model (GMM) based density estimation to generate the synthetic data. Then, it uses grow-and-prune DNN synthesis to generate inference models that are both accurate and computationally efficient. Section 3.3 discusses the MHDeep DNN synthesis process in detail. MHDeep generates DNN architectures that are efficient enough to be deployed on the edge devices such as smartphones or smartwatches. Section 3.4 discusses the inference process of the MHDeep DNNs for diagnosis and daily monitoring of mental health issues.

## 3.2 Data collection and preparation

We collected WMS data from a total of 74 adult participants at the Hackensack Meridian Health Carrier Clinic, Belle Mead, New Jersey. The participants were diagnosed by medical professionals at the clinic. The 74 participants comprised the following four categories: 25 healthy participants (no mental health disorder), 23 participants with bipolar disorder, 10 participants with major depressive disorder, and 16 participants with schizoaffective disorder. The experimental procedure for data collection and analysis was approved by the Institutional Review Board of Princeton University. The physiological signals of the participants were captured by a commercially-available Empatica E4 smartwatch [43] and a Samsung Galaxy S4 smartphone, as shown in Fig. 3. We summarize all the data types collected in this study in Table 1. The physiological signals are derived from WMSs embedded in the smartwatch. They include GSR that measures sympathetic nervous system arousal, IBI that indicates the heart rate, ST that provides skin temperature readings, and 3-axis accelerometer (Acc-W) that measures acceleration in the $X$, $Y$, and $Z$ directions. Ambient and motion information is captured using sensors in the smartphone that include ambient temperature (Temp), gravity (Grav), acceleration (Acc-P), and angular velocity (Vel). It is worth mentioning that the acceleration sensors in the smartphone and smartwatch have different sampling rates, and capture different motion information.

Before data collection, all participants are informed about the experiment and are asked to sign a consent form. The data collection setup consists of placing the Empatica E4 smartwatch on the wrist of the participant's non-dominant hand and placing the Samsung Galaxy S4 smartphone in the opposite front pocket. Data collection lasts around 1.5 hours, during which time the participant is allowed to freely move around in the room with their on-body devices. During this time, the smartwatch and smartphone continuously record and store the physiological signals and ambient/motion information. At the end of the data collection period, we remove the smartwatch from the patient's wrist and the smartphone from their pocket. We use the Empatica E4 Connect portal for smartwatch data retrieval. We use a private Android application to download the smartphone data streams. All of the recorded data are timestamped at the time of sampling.

TABLE 1
Data types collected in the MHDeep framework

| Data type | Sampling rate (Hz) | Data source |
| --- | --- | --- |
| Galvanic skin response ($\mu$S) | 4 | Smartwatch |
| Skin temperature ($^\circ C$) | 4 | Smartwatch |
| Inter-beat interval ($ms$) | 1 | Smartwatch |
| Acceleration ($x, y, z$) | 32 | Smartwatch |
| Ambient temperature (%) | 5 | Smartphone |
| Gravity ($x, y, z$) | 5 | Smartphone |
| Acceleration ($x, y, z$) | 5 | Smartphone |
| Angular velocity ($x, y, z$) | 5 | Smartphone |

Next, we preprocess the dataset for use in DNN training. We first synchronize the smartwatch and smartphone data streams for each participant. This is necessary since the WMS data streams may vary in their start times and frequencies. Then, we divide the data for each participant into 15-second windows. Each 15-second window of the combined smartwatch/smartphone data constitutes one data instance. There is no time overlap between data instances. To obtain each data instance, we flatten and concatenate the data within the same time window from both the smartwatch and smartphone. This results in a feature space of dimension 2325. The smartwatch (smartphone) contributes 1575 (750) features. All the smartphone sensors have a sampling rate of 5Hz. In addition, the smartwatch sensors include one data stream at 32Hz, two data streams at 4Hz, and one data stream at 1Hz.

For each classification task, since the number of individuals in each of the four categories is small, we created three different data partitions for evaluation. The data instances extracted from the individuals in each of the four groups (healthy, schizoaffective, depressive, and bipolar) were divided into three sets: training, validation, and test. To evaluate the models on different unseen patients, data instances included in the training, validation, and test sets came from different individuals, i.e., no individual contributed data to more than one of these sets. Among the healthy participants, for each of the three data partitions, data instances from 15 individuals (60% of the healthy participants) are selected for the training set, from 5 individuals (20% of the healthy participants) for the validation set, and from the remaining 5 individuals (20% of the healthy participants) for the test set. For individuals with bipolar disorder, the training, validation, and test sets contain data instances from 13, 5, and 5 participants, respectively. Among the participants who had major depressive disorder, data instances from 6 participants are selected for the training set and from 2 participants each for the validation and test sets. For individuals with schizoaffective disorder, the training, validation, and test sets include data instances from 10, 3, and 3 participants, respectively.

We create the final dataset for each binary classification task (healthy vs. the mental health disroder) by combining the training, validation, and test sets of the two classes involved in that task. We use SMOTE [44] to up-sample data instances from the minority class. Table 2 shows the number of instances for each of the classification tasks for all three data partitions.
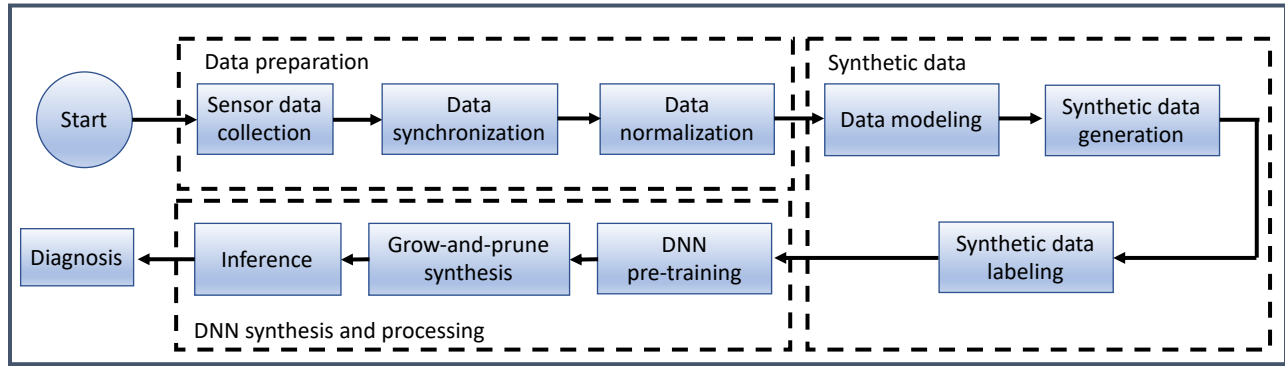
Fig. 2. Schematic diagram of the MHDeep framework.

TABLE 2
Details of various datasets (MDD: major depressive disorder).

| Classification task | Data partition | Training instances (#individuals) | Validation instances (#individuals) | Test set (#individuals) |
|---|---|---|---|---|
| Healthy vs. Bipolar | 1 | 14828 (28) | 3582 (10) | 3754 (10) |
| Healthy vs. MDD | 1 | 14828 (21) | 2414 (7) | 2515 (7) |
| Healthy vs. Schizo. | 1 | 14828 (25) | 2789 (8) | 3047 (8) |
| Healthy vs. Bipolar | 2 | 13330 (28) | 3922 (10) | 3582 (10) |
| Healthy vs. MDD | 2 | 13330 (21) | 3266 (7) | 2414 (7) |
| Healthy vs. Schizo. | 2 | 13330 (25) | 3773 (8) | 2789 (8) |
| Healthy vs. Bipolar | 3 | 12054 (28) | 4102 (10) | 3922 (10) |
| Healthy vs. MDD | 3 | 12054 (21) | 3088 (7) | 3266 (7) |
| Healthy vs. Schizo. | 3 | 12054 (25) | 3522 (8) | 3773 (8) |



Fig. 3. An Empatica E4 smartwatch (left) and Samsung Galaxy S4 smartphone (right) used in the data collection process.

### 3.3 MHDeep DNN synthesis

Fig. 4 shows the DNN architectures used in the MHDeep framework. The architectures receive the input data at the bottom and make their diagnostic decisions at the top. For the healthy vs. major depressive disorder and healthy vs. schizoaffective disorder binary classification tasks, the DNN architecture has four layers with a width of 256, 128, 128, and 2, respectively. For the healthy vs. bipolar disorder binary classification task, we use a DNN architecture with five layers with a width of 256, 128, 64, 32, and 2, respectively. We selected these DNN architectures by verifying the perfor-

mance of various DNNs (with different numbers of layers and number of neurons per layer) on the validation set and picking the best performing one. These architectures are initially fully-connected.
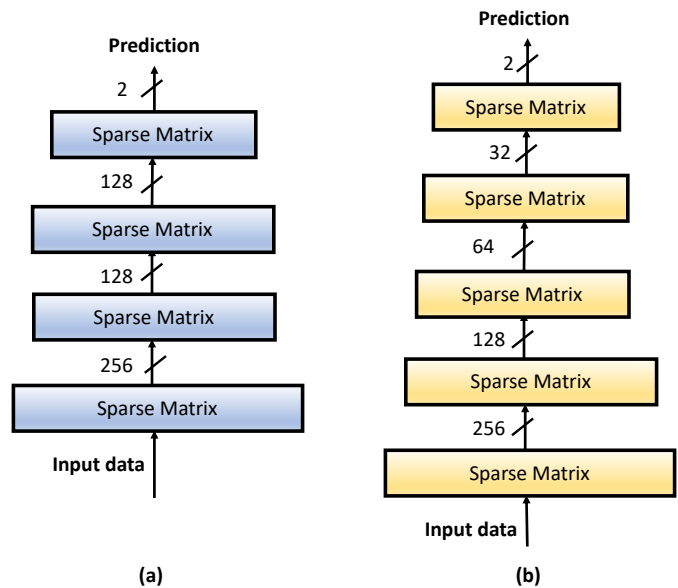


Fig. 4. Architecture of MHDeep DNNs for: (a) healthy vs. major depressive disorder and healthy vs. schizoaffective disorder, and (b) healthy vs. bipolar disorder.

We then subject the fully-connected DNNs to three sequential steps: (i) synthetic data generation to mimic the distribution of the real training data, (ii) pre-training of the

DNN architectures with the synthetic data, and (iii) grow-and-prune DNN synthesis to reduce the redundancy of the DNN model while improving its performance. Next, we discuss each step in more detail.

1) **Synthetic data generation**: In this step, we generate a synthetic dataset that mimics the probability distribution of the real training dataset. Fig. 5 illustrates the synthetic data generation process. This approach was first presented in [45] to alleviate the need for large datasets to train DNN architectures. We first use GMM to estimate the density of the training dataset. We optimize the number of mixtures in the GMM by monitoring the likelihood of validation data. We choose the number of components that maximizes the following criterion:

$$N^* = \arg\max_N \left(\text{GMM}_N(x).\text{score}(X_{validation})\right)$$

Finally, we train the optimal GMM model with $N^*$ mixtures on the combination of the training and validation datasets. By sampling this model, we are able to generate the synthetic data:

$$X^* = GMM_{N^*}(X_{total}).sample()$$

For our experiments, we generate 100,000 samples as synthetic data. The final step is labeling of the synthetic dataset. We use a traditional machine learning model for this purpose. We evaluate various models, e.g., the support vector machine and random forest models based on different splitting criteria (such as Gini index and entropy), and different depth limits on the decision trees, on the validation set. The model with the highest accuracy is used to label the synthetic data.

2) **DNN pre-training**: In this step, we use the labeled synthetic data to obtain a prior on the weights of the DNN architecture by pre-training them. The intuition behind this step is that pre-training the DNN provides a suitable inductive bias to the DNN. As a result, we can commence on the final training stage with a better weight initialization. Therefore, it alleviates the need for large training datasets.

3) **Grow-and-prune DNN synthesis**: MHDeep uses a grow-and-prune DNN synthesis paradigm to train the models. Algorithm 1 summarizes this process. It uses a mask-based approach. For each weight matrix, there is an associated binary mask of the same size that is used to disregard dormant connections in the architecture. It applies magnitude-based pruning and full growth to fully-connected DNNs iteratively. For magnitude-based pruning, a hyperparameter $\alpha$ is used to depict the pruning ratio. We prune a connection if and only if its weight is in the lowest $\alpha * 100$ percent of the weights in its associated layer. Finally, for the pruned connections, we set the weight and its binary mask both to 0. Since connection pruning is an iterative process, we retrain the network to recover its performance after each pruning iteration. In our experiments, after each architecture changing operation, we train the DNN

for 20 epochs. In addition, we set the number of iterations to 5.

---

**Algorithm 1** Grow-and-prune synthesis

---

**Input:** Pre-trained DNN architecture; iteration number $numIterations$; weight matrix $W \in R^{M \times N}$; mask matrix $Mask$ of the same dimension as the weight matrix; $\alpha$: pruning ratio

best-validation-acc = 0
**for all** layers in the DNN **do**
  $t = (\alpha \times MN)^{th}$ largest element in $|W|$
  **for all** $w_{ij}$ **do**
    **if** $|w_{ij}| < t$ **then**
      $Mask_{ij} = 0$
    **end if**
  **end for**
  $W = W \otimes Mask$
**end for**
Train DNN for given #epochs
validation-acc = evaluate DNN on validation set
**if** validation-acc >best-validation-acc **then**
  best-validation-acc = validation-acc
  Save the DNN
**end if**
**for** $numIterations$ **do**
  **for all** layers in the DNN **do**
    $Mask_{[1:M,1:N]} = 1$
  **end for**
  **for all** layers in the DNN **do**
    $t = (\alpha \times MN)^{th}$ largest element in $|W|$
    **for all** $w_{ij}$ **do**
      **if** $|w_{ij}| < t$ **then**
        $Mask_{ij} = 0$
      **end if**
    **end for**
    $W = W \otimes Mask$
  **end for**
  Train DNN for given #epochs
  validation-acc = evaluate DNN on validation set
  **if** validation-acc >best-validation-acc **then**
    best-validation-acc = validation-acc
    Save the DNN
  **end if**
**end for**
**Output:** Best architecture with the weight and mask matrices

---

### 3.4 MHDeep inference process

The trained DNN models can be used for diagnosis or daily monitoring of the mental state of the user based on collection of physiological signals and ambient information during the day. The collected data streams are processed based on the step explained in Section 3.2. We feed the processed data to the MHDeep DNN that predicts the mental health condition of the user. When the model predicts that the user is experiencing an episode of mental health disorder, this information can be sent to a physician for early treatment.
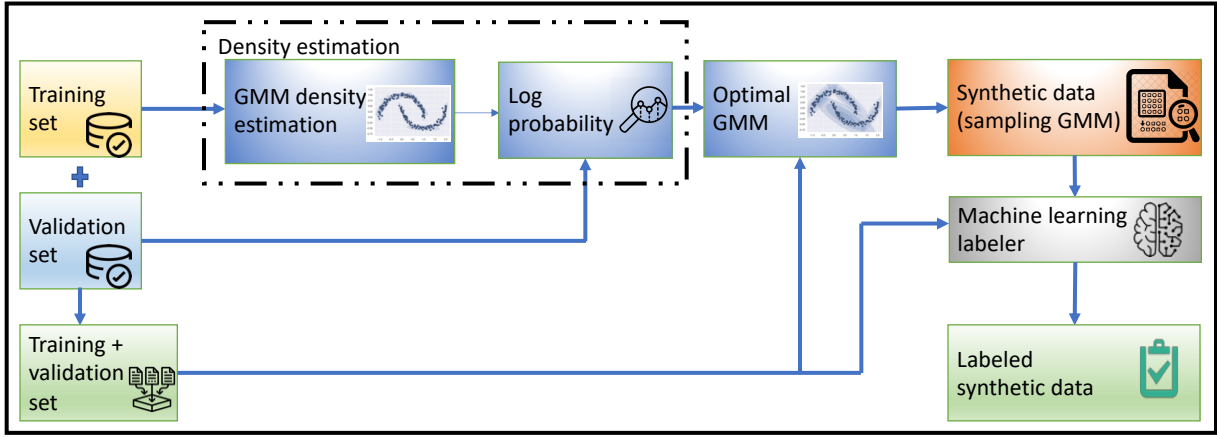
Fig. 5. Schematic diagram of the MHDeep synthetic data generation process.

## 4 IMPLEMENTATION DETAILS

We have implemented the data processing and preparation parts of the MHDeep framework in Python and the MHDeep DNN synthesis framework in PyTorch. We use the Nvidia Tesla P100 data center accelerator for DNN training. We use the cuDNN library to accelerate GPU processing. For training, we use a stochastic gradient descent (SGD) optimizer, with a learning rate of 5e-4 and a batch size of 256. We use 100,000 synthetic data instances to pre-train the network architecture. In the grow-and-prune synthesis phase, we train the network for 20 epochs each time the architecture changes. We use an SGD optimizer, with an initialized learning rate of 1e-4 that we halve in each succeeding iteration. We apply network-changing operations over five iterations.

## 5 EVALUATION

In this section, we analyze the performance of MHDeep DNN models for diagnosing three mental health disorders. This entails three binary classifications: (i) schizoaffective disorder vs. healthy individuals, (ii) major depressive disorder vs. healthy individuals, and (iii) bipolar disorder vs. healthy individuals. For each classification task, we use three different data partitions, each partition with data instances obtained from different individuals in the training, validation, and test sets.

The MHDeep DNN models are evaluated with four different metrics: test accuracy, false positive rate (FPR), false negative rate (FNR), and F1 score. Accuracy measures overall classification performance. It is simply the ratio of all the correct predictions on the test data instances and the total number of such instances. FPR and FNR measure how often healthy individuals are declared to have the corresponding mental health condition and *vice versa*, respectively.

First, we report the performance of the MHDeep DNN models in detecting each of the three mental health disorders at the data instance level. Next, we evaluate the accuracy of the models in detecting mental health disorders at the patient level.

## 5.1 MHDeep performance evaluation at the data instance level

We first analyze the performance of the three binary classifiers. We begin by training DNN models on features obtained from subsets of the eight data categories presented in Table 1. We analyze all the subsets of the eight data categories and report results for the top models. Since there are eight data categories, there are 256 subsets, with one being the null subset. We evaluated the remaining 255 subsets. This helps distinguish the impact of each data category and to find the most effective combination of categories for each classification task.

Table 3 shows the results for classification between healthy and schizoaffective data instances. The best data category subset, in this case, achieves an average test accuracy of 90.4%. We also report test accuracy, FPR, FNR, and F1 score for each of the three data partitions. The model reaches the highest test accuracy of 93.3% on the second data partition. Furthermore, for the healthy instances, the top model achieves a low average FPR of 6.5%, demonstrating its effectiveness in avoiding false alarms. For the schizoaffective instances, the model achieves an average FNR of 16.9%, indicating reasonable effectiveness in raising alarms when schizoaffective disorder does occur. We also report the number of parameters (#params) and floating-point operations (FLOPs) required for each model. We also compare #params and FLOPs of the models with those of the fully-connected baselines. As can be seen, using the grow-and-prune DNN synthesis approach enables us to reduce both #params and FLOPs, leading to a reduction in memory and computational requirements.

We present the results for classification between healthy and major depressive disorder instances in Table 4. The data category subset with the best performance achieves an average test accuracy of 87.3%. This model achieves the highest accuracy of 91.2% on the second data partition. It achieves an average FPR (FNR) of 6.8% (29.3%).

Table 5 presents the results for classification between healthy and bipolar disorder instances. In this case, the model trained on the best data category subset achieves an average test accuracy of 82.4%, with an FPR (FNR) of 16.7% (20.7%).

TABLE 3
Test accuracy, FPR, FNRs, and F1 score (all in %) for top data categories for classification between healthy and schizoaffective disorder data instances

| Data category | Data partition | #Params (compression) | FLOPs (compression) | Acc. | FPR | FNR | F1 Score |
|---|---|---|---|---|---|---|---|
| (Acc-P, Temp, Vel, Acc-W, GSR, IBI) | 1 | 275.0k (2.1×) | 549.5k (2.0×) | 91.2 | 11.4 | 5.8 | 89.5 |
| | 2 | 300.0k (1.9×) | 599.5k (1.9×) | 81.4 | 6.6 | 37.8 | 72.0 |
| | 3 | 275.0k (2.1×) | 549.5k (2.0×) | 87.3 | 2.4 | 34.1 | 84.2 |
| Average | | | | 86.6 | 6.8 | 25.9 | 81.9 |
| (Acc-P, Temp, Vel, Acc-W, GSR) | 1 | 200.0k (2.8×) | 399.5k (2.8×) | 90.5 | 13.0 | 4.4 | 89.3 |
| | 2 | 300.0k (1.9×) | 599.5k (1.8×) | 93.3 | 4.0 | 13.3 | 89.8 |
| | 3 | 250.0k (2.3×) | 499.5k (2.2×) | 87.5 | 2.5 | 32.9 | 77.9 |
| Average | | | | 90.4 | 6.5 | 16.9 | 85.7 |
| (Acc-P, Temp, Grav, Vel, Acc-W, GSR, IBI) | 1 | 300.0k (2.1×) | 599.5k (2.0×) | 88.3 | 14.4 | 7.7 | 86.8 |
| | 2 | 350.0k (1.8×) | 699.5k (1.8×) | 82.6 | 5.9 | 45.7 | 66.3 |
| | 3 | 300.0k (2.1×) | 599.5k (2.0×) | 88.7 | 3.9 | 26.4 | 81.1 |
| Average | | | | 86.5 | 8.1 | 26.6 | 78.1 |

TABLE 4
Test accuracy, FPR, FNRs, and F1 score (all in %) for top data categories for classification between healthy and major depressive disorder data instances

| Data category | Data partition | #Params (compression) | FLOPs (compression) | Acc. | FPR | FNR | F1 Score |
|---|---|---|---|---|---|---|---|
| (Temp, Grav, Vel, GSR) | 1 | 75.0k (2.7×) | 149.5k (2.4×) | 89.0 | 4.6 | 26.7 | 79.4 |
| | 2 | 120.0k (1.7×) | 239.5k (1.5×) | 90.5 | 4.5 | 21.7 | 82.7 |
| | 3 | 150.0k (1.3×) | 299.5k (1.2×) | 81.7 | 12.7 | 37.9 | 60.2 |
| Average | | | | 87.1 | 7.3 | 28.8 | 74.1 |
| (Acc-P, Grav, Vel, GSR) | 1 | 120.0k (2.0×) | 239.5k (1.8×) | 88.2 | 6.5 | 24.8 | 78.7 |
| | 2 | 145.0k (1.6×) | 289.5k (1.5×) | 91.2 | 1.9 | 25.7 | 83.0 |
| | 3 | 185.0k (1.3×) | 369.5k (1.2×) | 82.4 | 11.9 | 37.5 | 61.3 |
| Average | | | | 87.3 | 6.8 | 29.3 | 74.3 |

TABLE 5
Test accuracy, FPR, FNRs, and F1 score (all in %) for top data categories for classification between healthy and bipolar disorder data instances

| Data category | Data partition | #Params (compression) | FLOPs (compression) | Acc. | FPR | FNR | F1 Score |
|---|---|---|---|---|---|---|---|
| (Acc-P, Temp, Grav, Vel, Acc-W, IBI, ST) | 1 | 500.0k (1.2×) | 999.5k (1.2×) | 76.1 | 47.1 | 2.8 | 81.0 |
| | 2 | 480.0k (1.3×) | 959.5k (1.3×) | 81.4 | 1.0 | 34.3 | 78.9 |
| | 3 | 400.0k (1.6×) | 799.5k (1.5×) | 89.8 | 2.1 | 25.1 | 83.8 |
| Average | | | | 82.4 | 16.7 | 20.7 | 81.2 |
| (Temp, Grav, Vel, Acc-W, GSR, IBI) | 1 | 490.0k (1.2×) | 979.5k (1.1×) | 75.6 | 47.1 | 3.8 | 80.5 |
| | 2 | 450.0k (1.3×) | 899.5k (1.2×) | 81.6 | 0.9 | 34.2 | 79.0 |
| | 3 | 380.0k (1.5×) | 759.5k (1.5×) | 87.0 | 2.1 | 33.0 | 78.4 |
| Average | | | | 81.4 | 16.7 | 23.7 | 79.3 |

## 5.2 MHDeep performance evaluation at the patient level

Next, we show patient-level diagnostic test accuracy. We use the most accurate model from among the models discussed above for each classification task. Fig. 6 shows the results. In these graphs, we plot patient-level test accuracy vs. the duration of data needed for inference. Prediction is performed for each patient by simply taking the majority of the predicted labels for each 15s data instance in the given data duration. We step up the data duration size by 2 minutes each time. Thus, we add eight data instances in each 2-minute window. As we can see, the models reach 100% test accuracy after a certain point for distinguishing healthy individuals from those with schizoaffective and major depressive disorders. In addition, the best model for classification between healthy and bipolar disorder individuals reaches 90.0% patient-level accuracy. Table 6 shows the minimum data duration needed to reach the saturation accuracy. The durations are 40, 16, and 22 minutes for healthy vs. schizoaffective disorder, healthy vs. major depressive disorder, and healthy vs. bipolar disorder classifications, respectively.

## 6 DISCUSSION

MHDeep combines efficient neural networks with commercially available WMSs to diagnose various mental health disorders. Although several works address mental health problem detection using machine learning, to our knowledge, MHDeep is the only solution that focuses on an easy-to-use system that can monitor the daily mental health state of the
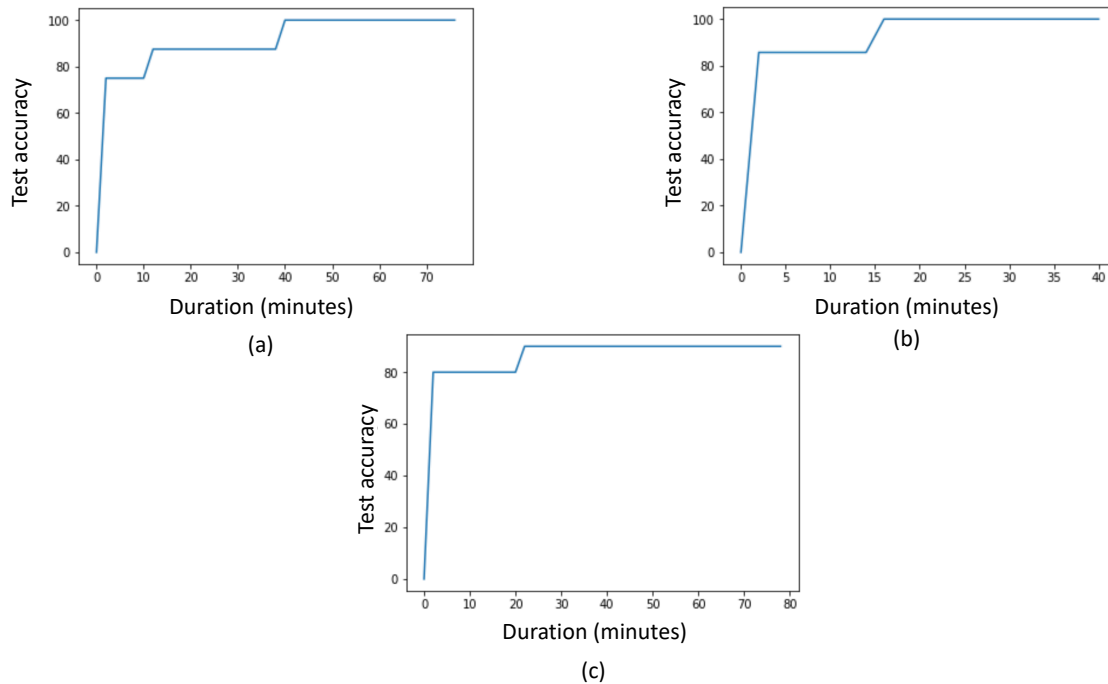
Fig. 6. Patient-level test accuracy vs. duration of data needed for classification between (a) healthy and schizoaffective disorder individuals, (b) healthy and major depressive disorder individuals, and (c) healthy and bipolar disorder individuals.

TABLE 6
Minimum inference data duration (in minutes) needed to reach saturation patient-level accuracy (in %) for each classification task

| Classification | Time | Saturation accuracy |
|---|---|---|
| Healthy vs. Schizoaffective disorder | 40 | 100 |
| Healthy vs. Major depressive disorder | 16 | 100 |
| Healthy vs. Bipolar disorder | 22 | 90.0 |

user. The diagnostic decisions can be sent to a health server from where medical professionals can access the information. This can enable them to quickly intervene during severe episodes of the disorder.

We demonstrated the diagnostic effectiveness of MHDeep for three different mental health disorders. However, this work can be generalized to other types of mental health disorders as well. We hope it will encourage researchers to start collecting WMS data from individuals across a diverse set of challenging diagnostic tasks. Bypassing the manual feature engineering step through the use of efficient DNNs enables easy scalability of this approach to many other disease domains. It can also be used to predict the progress of a disease based on longitudinal WMS data collected in the training stage.

Diagnosis of mental health disorders is often based on patient's self-report and answers to a questionnaire designed to detect each disorder. In the future, we can improve the performance of MHDeep by adding a specifically designed questionnaire to the data categories. In addition, adding features based on other WMSs such as blood pressure may also help enhance model performance.

# 7 CONCLUSION

In this article, we proposed a framework called MHDeep that combines data obtained from commercially available WMSs with the knowledge distillation power of DNNs for continuous and pervasive diagnosis of three main mental health disorders: schizoaffective, major depressive, and bipolar. MHDeep uses a synthetic data generation module to address the lack of large datasets. We trained the DNN models by using iterative network growth and pruning to learn both the weights and architecture during the training process. We evaluated MHDeep based on data collected from 74 individuals. It achieves patient-level accuracy of 100%, 100%, and 90.0%, using 40, 16, and 22 minutes of data collected in the inference stage, for classification between healthy and schizoaffective disorder individuals, healthy and major depressive disorder individuals, and healthy and bipolar disorder individuals, respectively. The MHDeep models were also shown to be computationally efficient. Thus, MHDeep can be employed for pervasive diagnosis and daily monitoring while offering high computational efficiency and accuracy.

# REFERENCES

[1] Z. Steel, C. Marnane, C. Iranpour, T. Chey, J. W. Jackson, V. Patel, and D. Silove, "The global prevalence of common mental disor-
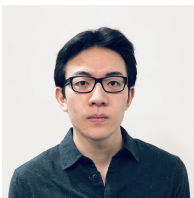
ders: A systematic review and meta-analysis 1980–2013," *Int. J. Epidemiology*, vol. 43, no. 2, pp. 476–493, 2014.

[2] G. V. Polanczyk, G. A. Salum, L. S. Sugaya, A. Caye, and L. A. Rohde, "Annual research review: A meta-analysis of the worldwide prevalence of mental disorders in children and adolescents," *J. Child Psychology and Psychiatry*, vol. 56, no. 3, pp. 345–365, 2015.

[3] D. B. Dwyer, P. Falkai, and N. Koutsouleris, "Machine learning approaches for clinical psychology and psychiatry," *Annual Review Clinical Psychology*, vol. 14, pp. 91–118, 2018.

[4] S. Hassantabar, N. Stefano, V. Ghanakota, A. Ferrari, G. N. Nicola, R. Bruno, I. R. Marino, K. Hamidouche, and N. K. Jha, "CovidDeep: SARS-CoV-2/COVID-19 test based on wearable medical sensors and efficient neural networks," *arXiv preprint arXiv:2007.10497*, 2020.

[5] S. Hassantabar, M. Ahmadi, and A. Sharifi, "Diagnosis and detection of infected tissue of COVID-19 patients based on lung X-ray image using convolutional neural network approaches," *Chaos, Solitons & Fractals*, p. 110170, 2020.

[6] R. Miotto, F. Wang, S. Wang, X. Jiang, and J. T. Dudley, "Deep learning for healthcare: Review, opportunities and challenges," *Briefings in Bioinformatics*, vol. 19, no. 6, pp. 1236–1246, 2018.

[7] H. Yin, B. Mukadam, X. Dai, and N. K. Jha, "DiabDeep: Pervasive diabetes diagnosis based on wearable medical sensors and efficient neural networks," *IEEE Trans. Emerging Topics in Computing*, 2019.

[8] A. O. Akmandor and N. K. Jha, "Smart health care: An edge-side computing perspective," *IEEE Consumer Electronics Magazine*, vol. 7, no. 1, pp. 29–37, 2017.

[9] H. Yin, A. O. Akmandor, A. Mosenia, and N. K. Jha, "Smart healthcare," *Foundations and Trends in Electronic Design Automation*, vol. 12, no. 4, pp. 401–466, 2018.

[10] H. G. Schnack, M. Nieuwenhuis, N. E. van Haren, L. Abramovic, T. W. Scheewe, R. M. Brouwer, H. E. H. Pol, and R. S. Kahn, "Can structural MRI aid in clinical classification? a machine learning study in two independent samples of patients with schizophrenia, bipolar disorder and healthy subjects," *Neuroimage*, vol. 84, pp. 299–306, 2014.

[11] U. R. Acharya, S. L. Oh, Y. Hagiwara, J. H. Tan, H. Adeli, and D. P. Subha, "Automated EEG-based screening of depression using deep convolutional neural network," *Computer Methods and Programs in Biomedicine*, vol. 161, pp. 103–113, 2018.

[12] J. Geraci, P. Wilansky, V. de Luca, A. Roy, J. L. Kennedy, and J. Strauss, "Applying deep neural networks to unstructured text notes in electronic medical records for phenotyping youth depression," *Evidence-based Mental Health*, vol. 20, no. 3, pp. 83–87, 2017.

[13] M. L. Birnbaum, R. Norel, A. Van Meter, A. F. Ali, E. Arenare, E. Eyigoz, C. Agurto, N. Germano, J. M. Kane, and G. A. Cecchi, "Identifying signals associated with psychiatric illness utilizing language and images posted to Facebook," *Nature Schizophrenia*, vol. 6, no. 1, pp. 1–10, 2020.

[14] X. Dai, H. Yin, and N. K. Jha, "NeST: A neural network synthesis tool based on a grow-and-prune paradigm," *IEEE Trans. Computers*, vol. 68, no. 10, pp. 1487–1497, 2019.

[15] S. Hassantabar, Z. Wang, and N. K. Jha, "SCANN: Synthesis of compact and accurate neural networks," *arXiv preprint arXiv:1904.09090*, 2019.

[16] J. F. Lehman, "The diagnostic and statistical manual of mental disorders," *Am. Psychiatric Assoc.*, 2000.

[17] P. W. Corrigan, "Mental health stigma as social attribution: Implications for research methods and attitude change," *Clinical Psychology: Science and Practice*, vol. 7, no. 1, pp. 48–67, 2000.

[18] J. E. Bordieri and D. E. Drehmer, "Hiring decisions for disabled workers: Looking at the cause," *J. Applied Social Psychology*, vol. 16, no. 3, pp. 197–208, 1986.

[19] P. W. Corrigan, D. Mittal, C. M. Reaves, T. F. Haynes, X. Han, S. Morris, and G. Sullivan, "Mental health stigma and primary health care decisions," *Psychiatry Research*, vol. 218, no. 1-2, pp. 35–38, 2014.

[20] C. Su, Z. Xu, J. Pathak, and F. Wang, "Deep learning in mental health outcome research: A scoping review," *Translational Psychiatry*, vol. 10, no. 1, pp. 1–26, 2020.

[21] D. Kuang and L. He, "Classification on ADHD with deep learning," in *Proc. Int. Conf. Cloud Computing and Big Data*, 2014, pp. 27–32.

[22] D. Kuang, X. Guo, X. An, Y. Zhao, and L. He, "Discrimination of ADHD based on fMRI data with deep belief network," in *Proc. Int. Conf. Intelligent Computing*, 2014, pp. 225–232.

[23] L.-L. Zeng, H. Wang, P. Hu, B. Yang, W. Pu, H. Shen, X. Chen, Z. Liu, H. Yin, Q. Tan *et al.*, "Multi-site diagnostic classification of schizophrenia using discriminant deep learning with functional connectivity MRI," *EBioMedicine*, vol. 30, pp. 74–85, 2018.

[24] W. H. Pinaya, A. Gadelha, O. M. Doyle, C. Noto, A. Zugman, Q. Cordeiro, A. P. Jackowski, R. A. Bressan, and J. R. Sato, "Using deep belief network modelling to characterize differences in brain morphometry in schizophrenia," *Scientific Reports*, vol. 6, p. 38897, 2016.

[25] G. Xiang-Fei and X. Jun-Hai, "Application of autoencoder in depression diagnosis," *DEStech Trans. Computer Science and Engineering*, no. csma, 2017.

[26] T. Pham, T. Tran, D. Phung, and S. Venkatesh, "Predicting healthcare trajectories from medical records: A deep learning approach," *J. Biomedical Informatics*, vol. 69, pp. 218–229, 2017.

[27] H. Lin, J. Jia, Q. Guo, Y. Xue, Q. Li, J. Huang, L. Cai, and L. Feng, "User-level psychological stress detection from social media using deep neural network," in *Proc. Int. Conf. Multimedia*, 2014, pp. 507–516.

[28] F. Sadeque, D. Xu, and S. Bethard, "UArizona at the CLEF eRisk 2017 pilot task: Linear and recurrent models for early depression detection," in *Proc. CEUR Workshop*, vol. 1866. NIH Public Access, 2017.

[29] G. Coppersmith, R. Leary, P. Crutchley, and A. Fine, "Natural language processing of social media as screening for suicide risk," *Biomedical Informatics Insights*, vol. 10, pp. 1–11, 2018.

[30] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNet v2: Inverted residuals and linear bottlenecks," in *Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.

[31] B. Wu, A. Wan, X. Yue, P. Jin, S. Zhao, N. Golmant, A. Gholaminejad, J. Gonzalez, and K. Keutzer, "Shift: A zero flop, zero parameter alternative to spatial convolutions," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2018, pp. 9127–9135.

[32] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," in *Proc. Int. Conf. Learning Representations*, 2016.

[33] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in *Proc. Int. Conf. Learning Representations*, 2017.

[34] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, "MnasNet: Platform-aware neural architecture search for mobile," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2019, pp. 2820–2828.

[35] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer, "FBNet: Hardware-aware efficient convnet design via differentiable neural architecture search," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2019, pp. 10 734–10 742.

[36] K. O. Stanley and R. Miikkulainen, "Evolving neural networks through augmenting topologies," *Evolutionary Computation*, vol. 10, no. 2, pp. 99–127, 2002.

[37] X. Dai, P. Zhang, B. Wu, H. Yin, F. Sun, Y. Wang, M. Dukhan, Y. Hu, Y. Wu, Y. Jia, P. Vajda, M. Uyttendaele, and N. K. Jha, "ChamNet: Towards efficient network design through platform-aware model adaptation," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2019.

[38] S. Hassantabar, X. Dai, and N. K. Jha, "STEERAGE: Synthesis of neural networks using architecture search and grow-and-prune methods," *arXiv preprint arXiv:1912.05831*, 2019.

[39] "Number of connected wearable devices worldwide from 2016 to 2022." [Online]. Available: https://www.statista.com/statistics/487291/global-connected-wearable-devices/

[40] M. M. Baig and H. Gholamhosseini, "Smart health monitoring systems: An overview of design and modeling," *J. Medical Systems*, vol. 37, no. 2, p. 9898, 2013.

[41] D. J. Malan, T. Fulford-Jones, M. Welsh, and S. Moulton, "CodeBlue: an ad hoc sensor network infrastructure for emergency medical care," in *Proc. Int. Workshop Wearable and Implantable Body Sensor Networks*, 2004.

[42] K. Wac, A. Van Halteren, and D. Konstantas, "QoS-predictions service: Infrastructural support for proactive QoS-and context-aware mobile services (position paper)," in *Proc. Int. Conf. Move to Meaningful Internet Systems*. Springer, 2006, pp. 1924–1933.

[43] ""Empatica E4 connect portal"." [Online]. Available: https://www.empatica.com/connect.

[44]  N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: synthetic minority over-sampling technique," *J. Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.

[45]  S. Hassantabar, P. Terway, and N. K. Jha, "TUTOR: Training neural networks using decision rules as model priors," *arXiv preprint arXiv:2010.05429*, 2020.

**Shayan Hassantabar** received his B.S. degree in Electrical Engineering, Digital Systems focus, from Sharif University of Technology, Iran. He also received his M.Math. degree in Computer Science from University of Waterloo, Canada, and his M.A. degree in Electrical and Computer Engineering from Princeton University. He is pursuing the Ph.D. degree in Electrical and Computer Engineering at Princeton University. His research interests include automated neural network architecture synthesis, neural network compression, and smart healthcare.

**Joe Zhang** received his B.S.E. from Princeton University in 2020. He is now pursuing a Ph.D. in Electrical Engineering at Stanford University, supported by a Stanford Graduate Fellowship. At Princeton, he was the recipient of the Shapiro Prize for Academic Excellence and the Hisashi Kobayashi Prize. He has worked on deep learning for healthcare applications and his current research at Stanford focuses on developing robust protocols and systems for blockchain technologies.

**Hongxu Yin** received his Ph.D. from Princeton University in 2020. He received his B.Eng. degree from Nanyang Technological University, Singapore, in 2015. He is now a Research Scientist with NVIDIA Research. He is a recipient of Princeton Yan Huo 94* Graduate Fellowship, Princeton Natural Sciences and Engineering Fellowship, Defense Science & Technology Agency gold medal, and Thomson Asia Pacific Holdings gold medal. His research focuses on efficient deep neural networks, data-free and hardware-guided model compression, and efficient inference for healthcare applications.

**Niraj K. Jha** received his B.Tech. degree in Electronics and Electrical Communication Engineering from Indian Institute of Technology, Kharagpur, India in 1981 and Ph.D. degree in Electrical Engineering from University of Illinois at Urbana-Champaign, IL in 1985. He has been a faculty member of the Department of Electrical Engineering, Princeton University, since 1987. He is a Fellow of IEEE and ACM, and was given the Distinguished Alumnus Award by I.I.T., Kharagpur. He has received the Princeton Graduate Mentoring Award.

He has served as the Editor-in-Chief of IEEE Transactions on VLSI Systems and an Associate Editor of several other journals. He has co-authored five widely used books. His research has won 20 best paper awards or nominations and 21 patents. His research interests include smart healthcare, cybersecurity, machine learning, and monolithic 3D IC design. He has given several keynote speeches in the areas of nanoelectronic design/test, smart healthcare, and cybersecurity.