

Deep Learning for Depression Recognition with Audiovisual Cues: A Review

Lang He^{a,b,*}, Mingyue Niu^{c,d}, Prayag Tiwari^{e,*}, Pekka Marttinen^{e,*}, Rui Su^{f,*}, Jiewei Jiang^{g,*}, Chenguang Guo^{h,*}, Hongyu Wang^{a,b}, Songtao Ding^{a,b}, Zhongmin Wang^{a,b}, Xiaoying Pan^{a,b}, Wei Dangⁱ

^aSchool of Computer Science and Technology, Xi'an University of Posts and Telecommunications, Xi'an Shaanxi 710121, China

^bShaanxi Key Laboratory of Network Data Analysis and Intelligent Processing, Xi'an University of Posts and Telecommunications, Xi'an Shaanxi 710121, China

^cNational Laboratory of Pattern Recognition (NLPR), Institute of Automatic Chinese Academy of Sciences (CASIA), Beijing 100190, China

^dSchool of Artificial Intelligence, University of Chinese Academy of Sciences (UCAS), Beijing 100049, China

^eDepartment of Computer Science, Aalto University, Espoo, Finland.

^fSchool of Foreign Languages, Northwest University, Xi'an Shaanxi, China

^gSchool of Electronic Engineering, Xi'an University of Posts and Telecommunications, Xi'an, China

^hSchool of Electronics and Information, Northwestern Polytechnical University, Xi'an Shaanxi, China

ⁱShaanxi Mental Health Center, Xi'an Shaanxi 710061, China

Abstract

With the acceleration of the pace of work and life, people have to face more and more pressure, which increases the possibility of suffering from depression. However, many patients may fail to get a timely diagnosis due to the serious imbalance in the doctor-patient ratio in the world. Promisingly, physiological and psychological studies have indicated some differences in speech and facial expression between patients with depression and healthy individuals. Consequently, to improve current medical care, many scholars have used deep learning to extract a representation of depression cues in audio and video for automatic depression detection. To sort out and summarize these works, this review introduces the databases and describes objective markers for automatic depression estimation (ADE). Furthermore, we review the deep learning methods for automatic depression detection to extract the representation of depression from audio and video. Finally, this paper discusses challenges and promising directions related to automatic diagnosing of depression using deep learning technologies.

Keywords: Affective computing, Depression, Deep learning, Automatic depression estimation, Review

1. Introduction

Depression is a type of mental disorder, which brings serious burden to individuals, families, and society. According to the World Health Organization (WHO), depression will be the most common mental disorder by 2030 [1]. In severe situations, depression leads to suicide [2, 3]. Besides, the report released by [3, 4] points out that approximately 50% of suicides are linked to depression. Currently, there is no unique and efficient clinical characterization of depression, which makes the diagnosis of depression time-consuming and subjective [5]. As gold-standard assessments or tools mainly depend on the subjective experience of clinicians, it is challenging to have a unified standard for diagnosing the severity of depression. The main diagnostic tools for severity, e.g., Hamilton Rating Scale for Depression (HAM-D) [6], rely on interviews conducted by clinicians or individuals themselves, yielding a score which summarizes the behavior of the patients. Diagnosing depression is complicated, depending not only on the educational background, cognitive ability, and honesty of the subject to describe the symptoms, but also on the experience and motivation of the clinicians. Comprehensive information and thorough clinical training are needed that to diagnose the severity of depression accurately [7]. Some biological markers, for instance, low serotonin levels [8, 9], neurotransmitter dysfunction [10, 11] and genetic abnormalities [12, 13], have been considered as indicators of depression,

*Corresponding author

Email addresses: langhe@xupt.edu.cn (Lang He), prayag.tiwari@aalto.fi (Prayag Tiwari), pekka.marttinen@aalto.fi (Pekka Marttinen), sabrina@nwu.edu.cn (Rui Su), jiangjw924@126.com (Jiewei Jiang), guochg@nwpu.edu.cn (Chenguang Guo)

however it is unclear which biomarker(s) is the most efficient as an indicator. Hence, in recent years, numerous Automatic Depression Estimation (ADE) systems have been introduced to automatically estimate the severity scale of depression based on audiovisual cues extracted with techniques developed in machine learning, speech recognition, and computer vision field [14, 15, 16, 17].

Designing a representative feature and its extraction for predicting the scale (i.e., severity) of depression is an important step in the deep learning architecture for ADE. ADE features can either be hand-crafted or based on deep learning models. Examples of widely used hand-crafted features include Local Binary Patterns (LBP) [18], Local Phase Quantization from Three Orthogonal Planes (LPQ-TOP) [19], Local Binary Patterns from Three Orthogonal Planes (LBP-TOP) [20], and others (e.g., Facial Action Units (FAUs), Landmarks, Head Poses, Gazes) [21]. However, since 2013, depression recognition challenges such as Audio-Visual Emotion Recognition Challenge (AVEC2013) [22] have recorded depression data by Human-Computer Interaction. Meanwhile, the fast development of deep learning has motivated many scientists to study Deep Learning (DL) approaches for depression recognition, which has resulted in a promising performance compared to the hand-crafted features. For deep learning approaches, a wide range of studies have adopted the Deep Convolutional Neural Network (DCNN) to extract multi-scale feature representations [23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33]. Fig. 1 shows this evolution of ADE according to methods and databases.

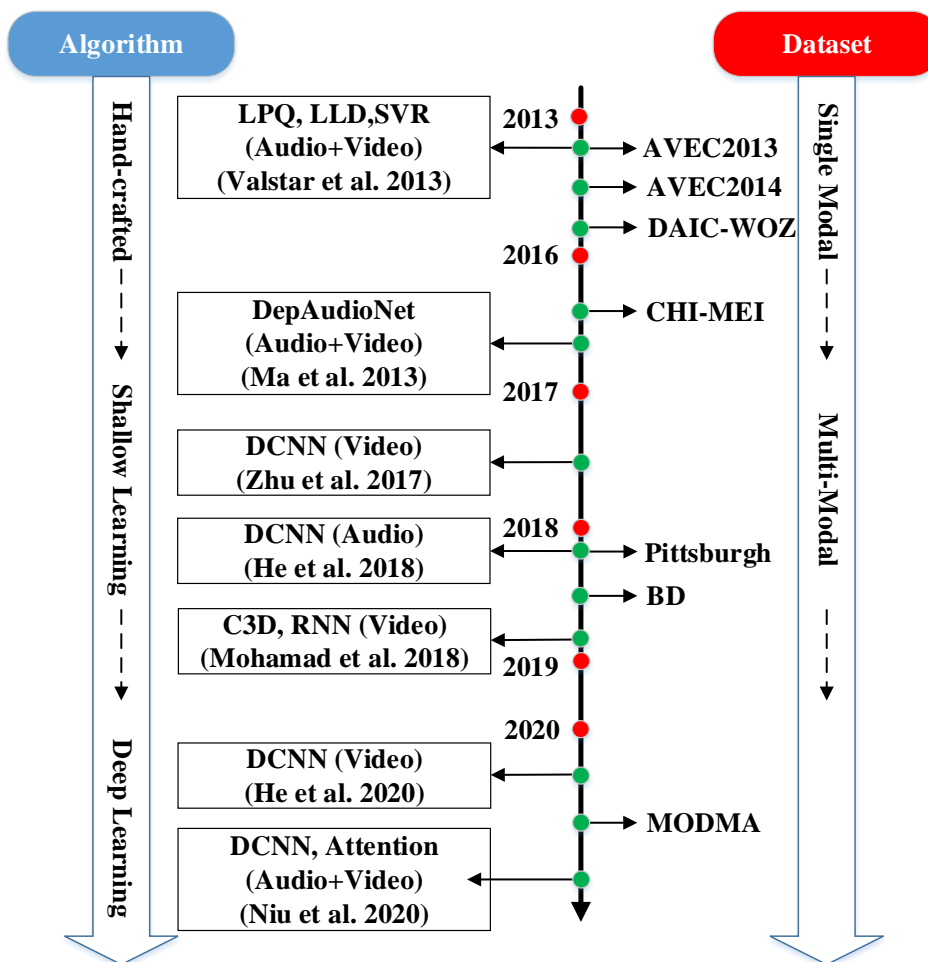


Figure 1: A brief chart to show the evolution of ADE with approaches and databases. From 2013 to 2021, the algorithm based on feature extraction had undergone three different stages from hand-crafted to shallow to deep learning. In the meantime the databases have evolved from single-modal (e.g., audio or video alone) to multi-modal (comprising multiple data types jointly).

1.1. Limitations

In recent years, some exhaustive reviews on depression recognition and analysis have been published based on audio [34], and visual cues [35]. These surveys provided a comprehensive scope for ADE. Yet, there remain two unexplored aspects. As existing reviews only focus either on audio or visual cues for estimating the depression scale, the combined use of audiovisual cues has not been adequately discussed. In addition, existing surveys only consider traditional approaches, and DL technology has not yet been covered in their analyses. Recently, DL technology has quickened the development and innovation of depression recognition based on audiovisual cues. As far as we know, an in-depth review on multi-modal audiovisual approaches for depression recognition is still missing. Our goal is to fill the gap in the existing extensive reviews by including the increasingly important deep multi-modal ADE approaches, based on audiovisual cues.

Despite the superior performance of deep learning, it still has some issues for ADE. First, training of deep learning methods (e.g. DCNN, Recurrent Neural Network (RNN), Convolutional 3D (C3D)) needs a lot of data and the existing public depression databases are limited and inadequate for the depression detection task. Second, individuals have different personal attributes, e.g., gender, age, race, educational background, etc. These personal attributes and traits have not been accounted for in existing deep learning models. For different individuals, these attributes allow considering depression from different perspectives. For instance, the education of an individual may indicate a broad interest in a lot of things, which may maybe preventive from depression. Third, combining hand-crafted and deep learning based features is a challenge that needs to be addressed to guarantee excellent performance in depression detection. Fourth, a fusion of multi-modal signals requires a solid theoretical foundation to mine the complementary representations from the different modalities.

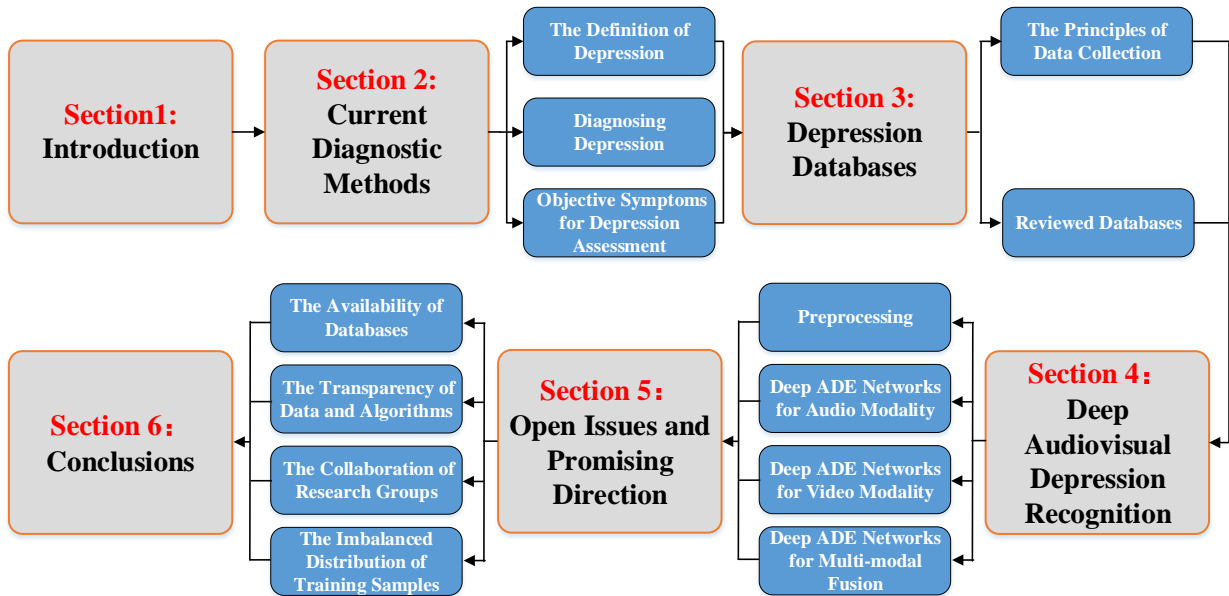


Figure 2: The paper structure through graphical illustration.

In this paper, we comprehensively review automatic depression detection methods based on deep neural networks, discuss the challenges, and point to future research directions. In the following, Section 2 provides the definition of depression and describes the objective markers for depression assessment. Section 3 introduces several multi-modal depression databases. Section 4 gives a detailed review of general deep ADE methods and presents several novel neural network architectures based on audiovisual cues. Additional issues are described in Section 5. Section 6 provides conclusions based on our exposition. In addition, to make a clear explanation, the structure of the paper is shown in Fig. 2.

2. Current Diagnostic Methods

To better understand the procedure of depression recognition based on audiovisual cues, the definition of depression is reviewed next, and then automatic methods to diagnose depression are surveyed.

2.1. The Definition of Depression

In 1980, Russell [36] proposes that emotional states can be represented as continuous numerical vectors in a two-dimensional space, called the Valence-Arousal (VA) space, see Fig. 3. The valence dimension refers to the two types of emotional states, i.e., positive and negative. The arousal dimension represents the intensity of emotion from sleepiness (or boredom) to high excitement. As shown in Fig. 3, depression is located in the third quadrant of the VA space [36], which corresponds to low-arousal and negative valence.

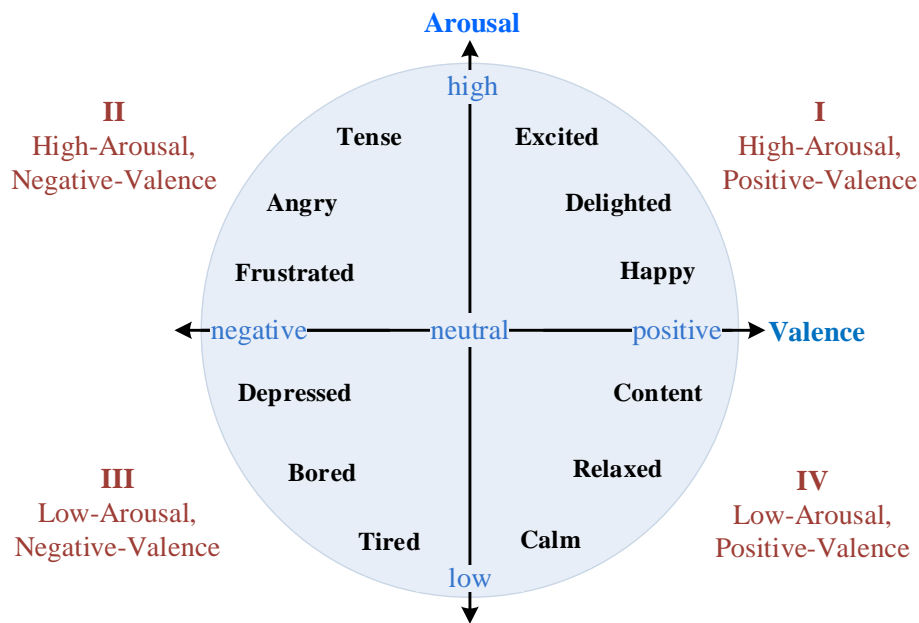


Figure 3: The two-dimension emotion space, which can be divided into four quadrants [37]. Every quadrant is associated with various emotions. For instance, high-arousal and positive-valence emotional states include excited, delighted, happy, etc. The depression disorder, which can be categorized in the third quadrant of the VA.

According to the definition from the Diagnostic and Statistical Manual of Mental Disorders (DSM) of the American Psychiatric Association (APA) [38], depression can be further divided as follows: Major Depressive Disorder (MDD), Persistent Depressive Disorder (Dysthymia), Disruptive Mood Dysregulation Disorder (DMDD), Premenstrual Dysphoric Disorder (PDD), Substance/Medication-Induced Depressive Disorder (S/M-IDD), Depressive Disorder Due to Another Medical Condition (DDAMC), and Other Specified Depressive Disorder (OSDD) or Unspecified Depressive Disorder (UDD). The DSM has provided the general criteria for classifying mental disorders by observed symptoms. When an individual has at least one of the following symptoms: 1) depressed mood most of the day and/or 2) markedly diminished interest or pleasure, in combination with at least four or more of the symptoms in Table 1 that have sustained for at least two weeks. In addition, the aforementioned symptoms are also expected to cause clinically significant distress or impairment in social, occupational, or other important areas of functioning. Nonetheless, these different types of depression related disorders manifest themselves in a similar way to a certain extent.

The question of how to diagnose depression has attracted attention of many researchers from different fields. However, the understanding of the pathogenesis of depression has not yet been unified and agreed upon. Nevertheless, pathogenesis is usually considered to be linked with a dysfunction of the cortical-limbic system, reducing its activity and connectivity [39, 40, 41, 42]. Nonetheless, it is believed that depression depends on the interaction between a genetic predisposition and environmental factors [43, 44]. For the effect of genetic predisposition, in [45], the authors

Table 1: Symptoms associated with depression [38].

Depressed Mood and/or Markedly diminished interest or pleasure

In combination with four of:

1. Significant weight loss when not dieting OR weight gain (e.g., a change of more than 5% of body weight in a month), decrease OR increase in appetite nearly every day. In children, a failure to make expected weight gains
 2. Insomnia OR hypersomnia nearly every day (inability to sleep OR excessive sleeping)
 3. Psychomotor agitation OR retardation nearly every day (observable by others, not merely subjective feelings of restlessness OR of being slowed down)
 4. Feelings of worthlessness OR excessive OR inappropriate guilt (which may be delusional) nearly every day (not merely self-reproach OR guilt about being sick)
 5. Diminished ability to think OR concentrate, OR indecisiveness, nearly every day (either by subjective account OR as observed by others)
 6. Fatigue OR loss of energy almost every day
 7. Recurrent thoughts of death (not just fear of dying), OR recurrent suicidal ideation without a specific plan, a suicide attempt OR a specific plan for committing suicide
-

found that monkeys suffered from depression may deprived of their mother. For the effect of environmental factors, in [46], Remi et al. found that for males, drinking too much in an adoptive family increases the risk of depression. For females, the death of an adoptive parent prior to adoptee age of 19 or the presence of an individual with a behavioural disorder in the adoptive family increases the risk of depression.

In addition, DSM has often been criticized that the boundaries between mental illnesses are not always properly defined. This resulted in subjective biases [47, 48, 49, 50, 51, 52]. There exist at the fewest 1497 unique profiles for depression [53]. In some cases of the same diagnosis, two depressed subjects may not share any identical symptom [54]. It is often considered that MDD can be referred to as *Clinical Depression*. As reported in the works of [55] and [56], the incremental economic loss of MDD has grown from 2005 to 2010 by 21.5% in the United States, while the economic loss is evaluated at 1% of the GDP.

2.2. Diagnosing Depression

It is difficult in primary care settings to assess the depression severity. Diagnosis of depression is often complicated with the chance of misidentification, the time-consuming nature, and the fact that not all depressed subjects directly show depressive manifestations (e.g., helplessness or hopelessness, etc.) [57, 58]. In addition, the combination of biological elements, family/environmental stressors, and personal vulnerabilities plays a vital role in affecting the onset of MDD [59].

Currently, the most frequently used assessment methods are interviews, e.g., HAMD [6], or self-assessments, e.g., the Beck Depression Index (BDI) (the first edition in 1961, and the recent version of 1996) [60]. A score is assigned by the assessment methodologies (HAMD and BDI) to each patient for characterizing their severity level by rating the 21 depression related symptoms. The main difference between HAMD and BDI is that HAMD requires a 20-30 min interview where the clinician fills the rated questionnaire, while BDI needs 5-10 min to complete the self-reported questionnaire. In addition, the two rating scales consider difference measures; the HAMD concentrates on neuro-vegetative symptoms (e.g., psychomotor retardation, weight, sleep, and fatigue, etc.), while the BDI focuses on self-assessments of negative self-evaluation symptoms. HAMD and BDI have been proven to obtain consistency when making a distinction between depressed from non-depressed patients [61, 62]. The HAMD tool was considered a gold standard for diagnosing the severity of depression. However, related research has exposed some issues [62, 63, 64]. Most importantly, some typical symptoms (i.e., insomnia, low mood, agitation, anxiety, and reduced weight)

Table 2: Commonly utilized depression rating scales [62].

Scale	Interview	Self Assessment	Number of items	Time to complete (Seconds)
HAMD [6]	√		17 or 21	20-30
BDI [60]		√	21	5-10
PHQ [76]		√	9	<5
QIDS [77]		√	16	5-10
MARSD [78]	√		10	20-30
Inventory of Depressive Symptomology (IDS) [79]	√		30	10-15
Zung Self-Report Depression Scale (Zung-SDS)[80]		√	20	5-10

related to the severity of depression are neglected by the HAMD. For every question in the HAMD questionnaire, the psychologists or clinicians should provide 3-5 possible responses to rate the severity of depression. A score in one of the ranges 0-2, 0-3, and 0-4 is assigned to indicate the severity of each symptom of depression. The summed score can be divided into five groups: (*Normal: the range from 0 to 7*), (*Mild: the range from 8 to 13*), (*Moderate: the range from 14 to 18*), (*Severe: the range from 19 to 22*) and (*Very Severe, ≥ 23*). While HAMD has covered numerous symptoms of depression, [65, 66] commented that only a part of the listed symptoms was useful for estimating the severity of depression. Nevertheless, a simple “symptom checklist” approach is considered insufficient to assess the ADE.

As mentioned above, the definition of depression from the clinical perspective may also depend on the scores provided by Self-report Scales and Inventories (SRSIs). The common assessment tools are BDI/BDI-II, PHQ-2/8/9 (Patient Health Questionnaire, 2, 8, or 9 is the number of questions, respectively), and Depression and Somatic Symptoms Scale (DSSS). To obtain a further comprehension of SRSI, BDI is introduced as follows. BDI is a commonly used assessment tool for SRSI of depression [67]. It consists of 21 questions, including cognitive, affective, and somatic symptoms, and several negative manifestations (e.g., self-evaluations and self-criticisms) related to depression. Every item of BDI/BDI-II, which is defined by multiple-choice options, is weighted by a numerical value (range: 0 – 3). The range of BDI scores is from 0 to 63 (*no or minimal depression: the range from 0 to 13*), (*Mild: the range from 14 to 19*), (*Moderate: the range from 20 to 28*), (*Severe: the range from 29 to 63*). Originally, BDI was not designed specifically for primary care usage, but its practical performance [68] shows that it is valid also for meetings in the primary care.

Though SRSIs have been widely used in various studies with the specificity and sensitivity reaching up to 80% to 90%, they nonetheless pose certain problems [69]. Specifically, the SRSI does not consider the clinical meaning of the observed symptoms and it allows individual variability when reporting different traits or characteristics, in contrast to a clinical interview [70]. Furthermore, SRSI cannot differentiate very well among different depression subtypes [61]. In addition, SRSI is susceptible to intentional or unintentional reporting bias [71]. Overall, despite the problem in providing an efficient diagnosis of depression [72, 73, 74], SRSIs have been widely adopted in various ways, for example in the primary health care, and research. Some researchers emphasized the cost-effectiveness SRSIs as a way for widespread screening practices to promote depression assessment [75].

Table 2 lists some depression scale ratings, e.g., HAMD, BDI, PHQ-9, Inventory of Depressive Symptomology (IDS), the 16-item Quick Inventory of Depressive Symptomology (QIDS), Zung Self-Report Depression Scale (Zung-SDS) and 10-item Montgomery–A Sberg Depression Rating Scale (MADRS), etc.

2.3. Objective Makers for Depression Assessment

It is generally considered that the representation of depression can be affected by various aspects [81, 82]. Observable behavioral signals were not accepted in the psychiatry field. However, several studies in these fields have still obtained popularity up to the present. Objective markers have been widely adopted in psychology, which can be utilized as an objective diagnostic tool in related fields (i.e., primary clinical settings, psychological institutions). It offered a robust assessment tool for assisting the clinicians in diagnosing the severity scale effectively and offered

later feedback and valuable advice for susceptible individuals. With the development of wearable devices, an interactive virtual tool is designed to deploy on the smartphone platform to help diagnose the depression subjects or the suspected people [17]. Therefore, there is an urgent requirement for designing novel assessment tools, such as developing diagnostic tools to investigate new markers. Previous studies on objective physiological, biological, and behavioral markers have improved the efficiency of psychiatric diagnosis and, they have the latent capacity to decrease the socio-economic costs caused by depression [54, 83].

In the early works of Emil Kraepelin [84], who was recognized as the father of modern psychiatry, he defined depressed voice as “patients speak in a low voice, slowly, hesitatingly, monotonously, sometimes stuttering, whispering, try several times before they bring out a word, become mute in the middle of a sentence”. In [34], the speech was considered a key objective maker for depression analysis, covering a wide range of features (i.e., prosodic, source, acoustic, and vocal tract dynamic).

In addition, patterns around facial regions were also significant for depression estimation. Hands and body posture are included in certain patterns related to depression assessment. Visual cues are essential for depression detection. It is considered that pupil dilation has a relation with depression. In [85], the authors consider that faster pupillary represents positive by healthy controls. Depressed subjects represent slower pupil dilation responses in certain conditions [86, 87, 88, 89, 90, 91]. In [92], the authors found that pupil bias and diameter are also important for assessing depression symptoms. In addition, the factors of facial expressions (e.g., anger, sadness, joy, surprise, disgust, fear, etc.) were regarded as a discriminative cue for depression detection. Suppose an individual is diagnosed to have depression symptoms. In that case, they will indicate low expressibility when showing the facial expressions [81, 93, 94, 95, 17, 96, 97, 89]. The features consist of reduced eye contact [98], gaze direction [99, 100, 101, 97], eyelid activity [102], iris movement [103], and eye openings/blinking [89, 103, 104]. Moreover, Eye movement and blinking were also considered a discriminative feature for predicting depression [102]. Furthermore, the duration of spontaneous smiles [99, 101], the intensity of smile [99, 97, 101], mouth animation [104], the lack of smiles [98] were also considered that contain valuable patterns for depression detection.

Action Units (AU) are the basic actions of muscle groups or individual muscles, which was originally proposed by Ekman et al. [105], and then adopted by Cohn et al. [14] to analyze the depression state. Meanwhile, a new AU-based method named Region Units (RU) was proposed [106]. Region Units (RU) are used to represent the regions of the face that enclose AUs. Various works have adopted AU to estimate the scale of depression and obtained promising performances [93, 94, 95, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119]. Later, it was found [102] that head pose and movement also contained discriminative patterns for assessing the severity of depression [93, 94, 95, 17, 99, 89, 110, 96, 97, 101, 103, 104, 120, 121, 122, 123]. Specifically, there were 46 points used to train a 3D face model using Active Appearance Models (AAM) to extract the head pose and movement features. Additionally, facial animation and the variability of motor were also utilized for depression detection [93, 95, 17, 104]. Body posture (i.e., the upper part of the body, the lower part of the body, and hands) are also very significant feature for detection [121, 122, 123, 124, 125]. In [100], the researchers considered that foot tapping and self-adapters also influence depression detection. In addition, the activity from facial muscles, skin electrical reactions, and peripheral blood pressure also caused involuntary changes, which often reflect the common and persistent negative thoughts and sad emotions of depression. In [126], they found that electroencephalographic recordings may have certain patterns related to depression. In [127] and [128], Functional Near-Infrared Spectroscopy (fNIRS) is also considered as a tool to assist in the depression diagnosis task [129].

In addition, several works have indicted that depression can cause the change of neurophysiological and neurocognitive abnormalities, which are demonstrated from an individual communication via facial gesturing, vocal articulation and so on [130, 35]. Therefore, we will focus on audiovisual signals for ADE in this review.

3. Depression Databases

Deep depression recognition requires sufficient data for the training of a discriminative model. Due to sensitivity of depression, data collection is challenging. Accordingly, different research groups have attempted to collect their own database to study assessment tools for depression estimation. Here, we introduce existing databases that have been widely adopted in the reviewed works for depression detection. Moreover, we also cover other privately released databases. Table 3 summarizes the mentioned databases, comprising the number of subjects, the annotation scores, ground truth, availability, additional details. Fig. 4 shows some example images from the AVEC2014 database.



Figure 4: Randomly selected images from video clips along with their BDI-II and depression severity scores on the AVEC2014 database [141]. To protect the privacy of participants, the images were blurred, and the regions of the eye were occluded. From the different rows, one can see that the severity of depression is enhanced from none to severe in the images.

3.1. The Principles of Data Collection

Collecting depression data requires recruiting a number of participants from hospitals or psychological clinics, which is the most challenging part of depression research. As reviewed in the existing studies, the depressed subjects or health controls are assessed by complying with DSM-IV [38] standard [85, 86, 88, 14, 110, 123] and/or HAMD [14, 111, 123, 125]. Also, Mini International Neuropsychiatric Interview (MINI) [147] was adopted to diagnose the severity of depression, and QIDS-SR was adopted for defining it. BDI has been applied extensively for predicting the scale of depression [85]. In other cases, several standards, PHQ-9 [93, 94, 95, 17, 99, 100] and BDI-II [22, 141] are aimed to assess the symptoms related to depression. Furthermore, other recruitment approach (e.g., flyers, posters, social networks, personal contacts, and mailing lists) have been utilized in several studies.

In order to obtain valuable patterns for depression prediction, the experimental environment should be finely designed. In general cases, some agreements were signed prior to the experiment. If the data collection occurs in hospitals, then some devices should be arranged first (e.g., cameras, microphones, and sensors), and the details of variable to record should be planned. Next, the details of the participants should be gathered (e.g., anamnesis record sheet and cognitive abilities). For instance, in [148], patients had to satisfy the following standards: 1) diagnosis with MDD or other mental disorder; 2) ability to understand and satisfy the protocol requirements; 3) no other clinical background that could disturb the results (e.g., delirium, dementia, amnesic or other symptoms); 4) no bipolar disorder symptoms assessed; 5) not the criteria of DSM-IV satisfied in the past three months; 6) ability to understand American English. In addition the the healthy control group had to comply with the same principles. Specifically, the healthy individuals should not have had any symptoms related to depression in the past year. All audio samples were collected in the same year and in an identical environment (i.e., room and other experimental setups).

As mentioned above, the collection environment or setup is essential for recording the data. In some scenarios, emotion elicitation was adopted to to produce a particular emotional response in the participants, which is different between healthy controls and depressed subjects [34]. Moreover, interviews have also been used to discover symptoms of depression, and some spontaneous emotional patterns closely related to depression [106] have also been identified during the interviews. Overall, interviews have been conducted by clinicians, psychologists, psychiatrists, and virtual human interviewers and also piloted by a computer to generate several data samples.

Regarding modality, speech and video samples [14, 94, 113, 17, 96, 97, 99, 100, 101, 104, 121, 122, 123, 149,

150, 16, 151, 152, 153], physiological signals [100, 154, 154, 155, 156, 157, 158, 159, 160], and text [100, 108] have been employed to improve the performance of depression assessment. However, different modalities were determined by the device used in the data collection stage. For an audio clip, a computer or laptop has been used to record the data samples (i.e., AVEC2013 [22], AVEC2014 [141], AVEC2016 [116]). For the video modality, the number of cameras and other attributes (i.e., colors and angles, etc.) have been used; for instance, the face and the whole body have been recorded separately with multiple cameras from different angles [14]. Also, thermal images based on eye temperature have been adopted to predict the severity of depression [161]. Microsoft Kinect has also been used to record for the upper body of a participant [93, 99]. The distance between participants and is approximately one meter. A portable three-electrode EEG device has been designed to collect the electroencephalography data [159]. Similarly, identical to emotion elicitation, the detailed setup alters across different works.

3.2. Reviewed Databases

1 – DementiaBank database [131] contains 226 persons from 2 longitudinal clinical-pathological studies in 4 years in USA between 1994 and 1998, and can be shared for the purpose of research. Video samples were recorded in this database. In the DementiaBank, all persons were diagnosed to suffer Alzheimer’s Disease (AD) at some stage. Besides, HAMD depression scores were labeled for the part of the participants from the subset of the database.

2 – The database [132] contains 43 depressed subjects (23 females and 20 males) from the hospital in 1998. Speech samples were recorded before the psychiatric exploration. This database is only accessed by themselves. In addition, most of the patients (79.1%) represented the onset of improvement in the 12 days, and they had nothing to do with the background of depression severity.

3 – The database [133] comprises of 115 peoples in USA. Audio modality is collected from different scenarios. This database is only accessed for lab usage by themselves. The male consisted of seventeen dysthymic patients, twenty-one major depressed patients, as well as ten control subjects. Of these participants, twenty-two high-risk suicidal patients and twenty-four control subjects. F0 and formants features were extracted from audio samples.

4 – The database [134] consists of a 12-week double-blind treatment trial in USA. Speech features were extracted in this database. 10 depressed subjects and 19 health controls were recruited. Depressed subjects were segmented into two groups by complying with the ratings, i.e., retarded or agitated. This database is also for laboratory use only.

5 – This database [135] was collected from at the Medical College of Georgia (MCG) at Augusta University in United States. It comprises 33 subjects in this task. Audio samples are collected in this database. All subjects should speak American English fluently and have not abused a significant substance. For each subject, their speech was recorded when reading a short story. Each speech recording fell to 65 separate sentences over the entire session. This database is also for laboratory use only.

6 – In this database [136], three different patient groups (i.e., depressed subjects, high-risk suicidal patients, and remitted patients) were included. Speech data samples were recorded of every subject. In addition, Power Spectral Densities (PSD’s) features were extracted from the speech data samples. This database is also for laboratory use only.

7 – This database [14] contains 57 participants (20 males and 37 females) from a clinical trial to diagnose depression symptoms in USA. All of them should met DSM-IV criteria of MDD during the clinical interview step. The severity symptom of depression is assessed on four occasions at approximately seven weeks intervals by a clinical interviewer. HRSD [162] is adopted to restrict the standard of interviews. Audio and video samples were recorded of the participants in this task.

8 – ORI database [137] comprises of eight participants during family interaction activities. Three interactions (Problem-solving Interaction (PSI), Event Planning Interaction (EPI), and Family Consensus Interaction (FCI)) were performed in the sessions. Every adolescent was asked to record a one-hour video recording during the sessions. The age of adolescents ranged from 12 to 19. The eight subjects fell into two groups, i.e., depressed subjects (4) and non-depressed groups (4). Of all of the participants were selected with white, and nobody wore glasses in the sessions. This database is not opened for public use.

9 – ORYGEN database [139] collects from ORYGEN Youth Health Research Center in Australia. The video and audio data samples are recorded from the discussions between parents and their children (the age ranging from 12 to 13 years) by two different family interactions (i.e., EPI, PSI). The mentioned video samples were recorded in two steps. In the first step (T1), 191 non-depressed subjects (94 females and 97 males) were recruited. Two years later, 15 (6 males & 9 females) of 191 suffered from MDD and three participants (1 male & 2 females) suffered from Other Mood Disorders (OMD) in the second step (T2). This database is not opened for public use.

10 – BlackDog database [138] is collected from an organization named BlackDog institute, focusing on a clinical study in Sydney, Australia. 80 participants (the age ranging from 21 to 75) participated. To ensure the availability of experiments, all participants had to comply with the criteria of DSM-IV. Speech data are recorded during the conversation between the interviewer and participants. The clinical interaction was performed by asking specific questions (eight groups), in which participants were required to describe events stimulated by specific emotions.

11 – This database [140] comprises 165 adults (104 females and 61 males) with major depression disorders between November 2006 and August 2007 in the United States. Speech samples were recorded with automated telephone devices in this database. All participants should meet the following criteria: 1) the age between 18 and 65; 2) not taking psychotropic medications; 3) HAMD of 22 or greater; 4) Symptoms of MDD has been diagnosis with DSM-IV, and have been lasted for a month. This database is also used in their own study.

12 – AVEC2013 database [22] refers to a selection from the audiovisual depressive corpus, covering 340 videos from 292 people by performing a human-computer interaction. 31.5 years (18 and 63 years) was the average age of the participants. BDI-II was adopted for labeling every audio and video segment. In this database, the organizer only provided a total of 150 audio and video clips, falling to three equivalent partitions (training, development, and test set). Different from the above mentioned databases, AVEC2013 is open for researchers to design the ADE systems.

13 – AVEC2014 corpus [141] was divided from the AVEC2013. The only difference was that the AVEC2014 database contained two tasks, i.e., Freeform and Northwind. Thus, every partition covered 100 data samples, respectively. Therefore, AVEC2014 contains 300 data samples in total. BDI-II was used for labeling every audio and video clip.

14 – Crisis Text Line. A web python development tool, back-end Flask, was adopted to store the data using MySQL database [142]. A single URL of the website was exploited to call with Flask routing between the front-end and back-end through AJAX. The interaction occurred between Flask and MySQL database using a Python class. Furthermore, the Python class could process the generated data accessing into and out of the database. Text features are extracted in this database.

15 – DAIC database [100] was collected based on a semi-structured clinical interactions in USA. Four types of interviews were conducted, i.e., Face-to-Face, Teleconference, Wizard-of-Oz, and Automated. The database contains 189 sessions of interactions, and consists of the audiovisual cues, as well as physiological data (e.g., galvanic skin response (GSR), electrocardiogram (ECG), and respiration). In addition, text modality was also collected during the interactions. Different verbal and non-verbal features were used to annotate the corpus. DAIC is the same as AVEC2013 and AVEC2014 to open access for the researchers.

16 – Rochester database [89] consists of 32 participants under different conditions. The normal group involved 27 participants (16 females and 11 males) in this database. The age of this group ranged from 19 to 33. Also, 5 depressed patients (2 severe and 3 moderate) and a healthy control group of five participants were covered.

17 – In the CHI-MEI database [143], six discrete videos (i.e., disgust, fear, sadness, surprise, anger, and happiness) were adopted to arouse the subjects to express their expressions based on their facial region and the responses speech of them. The CHI-MEI speech database was collected from the speech response of the subjects by a clinician in CHI-MEI Medical Center, Taiwan. The audio and video data were collected in this dataset. In total, 15 BDs, 15 UDAs and 15 healthy controls are recruited in CHI-MEI. In addition, the participants had to complete a baseline recording before data collection. Afterward, the participants watched six emotional videos.

18 – Pittsburgh database [144] comprises of 57 (34 females, 23 males) depressed participants from a clinical treatment for depression. The age ranged from 19 to 65 years (mean=39.65). All the participants had to reach DSM-IV criteria for MDD. Severity of MDD was assessed at 1, 7, 13, and 21 weeks by 10 random clinical interviewers. This database is also open for public usage.

19 – BD database [145] consists of 46 patients and 49 healthy controls of the mental health service of a hospital. To gather the sociodemographic and clinical patterns, all patients should perform semi-structured interviews by the SKIP-TURK. Young Mania Rating Scale (YMRS) and MADRS were employed to estimate the depressive and manic features in the next following days (0, 3, 7, 14, 28), and then altered in the third month. During this step, audiovisual samples were recorded. Accordingly, every video session was annotated by bipolar mania/depression ratings. This database is used as the challenge data in AVEC2018.

20 – MODMA database [146] was collected from audio and EEG signals for mental disorder analysis in China. Experienced psychiatrists rigorously recruited all the participants from hospitals. The EEG database contains data samples recorded with traditional 128-electrodes mounted elastic cap as well as a new wearable 3-electrode EEG

recorder for pervasive usage. The 128-electrodes EEG signals were recorded with resting-state and under-stimulation from 53 subjects, while 3-electrode EEG signals were recorded with resting-state from 55 subjects. Specific to the audio data, the samples were collected from 52 subjects by allowing participants to be interviewed, read stories, and watch the emotional pictures.

According to the mentioned databases, the following discussions are made:

- From the perspective of openness, most of the databases were only used for their own research and not released publicly for depression recognition study. Only few databases were released publicly for depression recognition, i.e., AVEC2013 and AVEC2014 ¹, DAIC-WOZ ², Pittsburgh dataset ³, and MODMA dataset ⁴. The rest of the databases may be available for the researchers in some scenarios.
- Most of the databases were collected by the region of the US and EU. There is only one database available for researchers in China, which is MODMA.
- From modalities, most of the database involved one or/and more (e.g., audio, video, physiological signals, text).
- In terms of the number of subjects, all the databases consisted of limited data samples, which is explained as depression is a mental disorder and also kept as a secret by the depressed subjects.

4. Deep Audiovisual Depression Recognition

This section presents the common procedures adopted in ADE, i.e., preprocessing, deep feature extraction, and classification/regression. For instance, we describe raw audio and other audio features that are used broadly as input for depression recognition, as well as other data modalities (e.g., video, text from transcripts, and physiological signals). In the following, the works introduced in the literature are separated into three groups: 1) deep ADE networks for audio modality; 2) deep ADE networks for static images, and 3) deep ADE networks for image sequences. In addition, different network types are also introduced for the mentioned groups along with discussion.

4.1. Preprocessing

Both conventional and end-to-end schemes require some pre-processing steps before the actual depression recognition and analysis.

For the audio data, the sample rate is usually processed to 16 kHz or others (e.g., AVEC2013). Meantime, to generate the spectrograms of audio data, Discrete Fourier Transform (DFT) method is adopted for conducting a Time-Frequency (TF) characterization for audio signals. To select DFT parameters, a Hanning window (23ms with 50% overlap) is used [33]. In addition, to extract efficient hand-crafted features, the length of low-level descriptor (LLD) is also considered in depression recognition studies. In [33], they tried LLDs of different lengths and suggested that 20s is sufficient to obtain a good performance. In [163], the authors sampled the waveforms at 8KHZ and generated the 129-dimensional normalized amplitude spectrogram using a short-time Fourier transform with 32 ms Hamming window and 16 ms frame shift for AVEC2013 and AVEC2014 databases.

For video data, image normalization, face detection, and alignment between the adjacent frames, are commonly used preprocessing techniques. Viola and Jones proposed a general algorithm for face detection [164]. Besides, the OpenFace toolkit provides a free tool for face detection and alignment in many applications [165]. The computer expression recognition toolbox is used in numerous fields, but it is not free to use at the moment [166]. Moreover, a comprehensive facial pre-processing tool can be found on the web ⁵. In addition, video data is pre-processed with different types for depression recognition, i.e., frame-level images, and image sub-sequences, and image sequences.

¹<http://avec2013-db.sspnet.eu/>

²<http://dcapswoz.ict.usc.edu/>

³<http://www.pitt.edu/emotion/depression.html>

⁴<http://modma.lzu.edu.cn/data/index/>

⁵<http://nordicapis.com/20-emotion-recognition-apis-that-willleave-you-impressed-and-concerned/>

4.2. Deep ADE Networks for Audio Modality

In the databases mentioned above, the extraction of hand-crafted features stays dominant in audio-based ADE. Next we describe hand-crafted feature extraction from audio cues for ADE.

Since 1998, a range of feature representation methods have been proposed to estimate the severity of depression. Here we list only some classical (shallow) methods for depression recognition, and then focus primarily on the deep automatic depression recognition methods. In 1998, the pause duration of speech was tightly related to the HAMD scores for 60% of patients [132]. In 2001, Alpert et al. [134] found that there are differences between healthy controls and depressed individuals. Cannizzaro et al. found an important relationship between reduced speak rate and HAMD score in 2004 [167]. In addition, they found that different acoustic features also could impact the representation of depression (e.g., percent pause time, speaking rate, and pitch variation). It is noteworthy that variations of speaking rate and pitch were considered an important representation for depression analysis. In 2008, Moore et al. [168] studied the combination of a wide range of features, e.g., prosodic, voice quality, spectral, and glottal. They obtained comparable performance for classifying the absence/presence of depression [168]. Many LLD indicators (e.g., prosodic, source, formant, and spectral) have been identified as efficient predictors of depression. For an in-depth review of speech-based depression recognition, please refer to the article [34]. As revealed from the review work, one can note that hand-crafted features have achieved promising performance for depression prediction. However, some issues remain; for instance, manual work and expert knowledge are significant for feature selection, which wastes labor resources. Furthermore, representations learned via DL have exhibited excellent performance compared to hand-crafted in multiple disciplines, and ADE is not an exception.

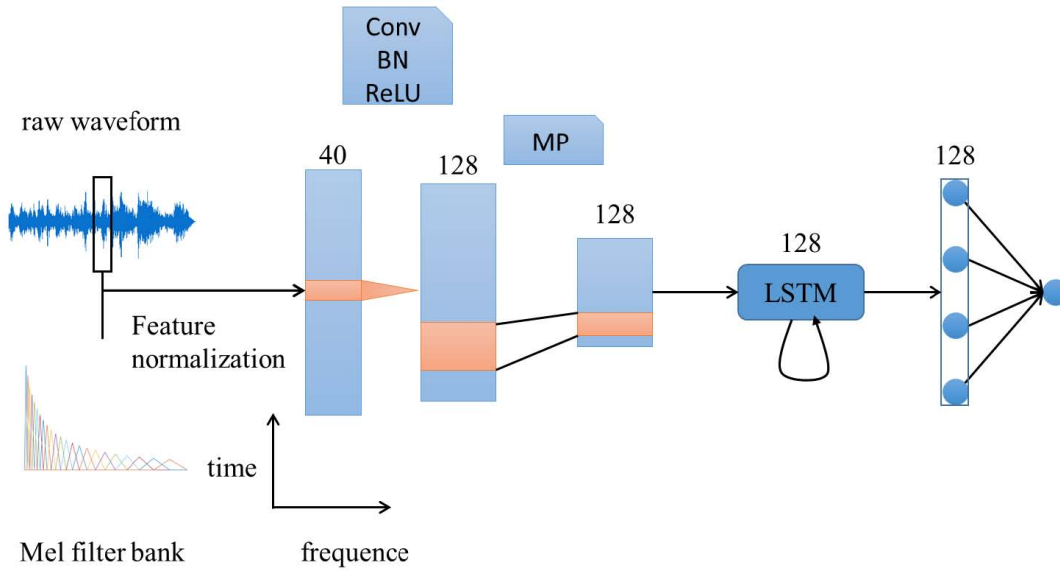


Figure 5: The pipeline of DepAudioNet framework from [23]. In the framework, CNN and LSTM are combined to model different scale features for audio-based depression recognition. The DCNN can model the high-level patterns of the raw waveforms. The LSTM can learn the combination of short-term and long-term representations from Mel-scale filter bank features. Mel-scale filter bank feature is used as a LLD to represent the characteristics from the vocals. Conv is the convolutional operation, BN represents the batch normalization, ReLU is the rectified linear unit operation, and MP is the multi-layer perceptron.

In 2016, [23] proposed a new model based on deep learning, i.e., DepAudioNet, to mine the depression representation from vocal cues, adopting LSTM and DCNN to encode a discriminative audio representation for depression recognition (see Fig. 5). DCNN can model the spatial feature representations from the raw waveforms, and LSTM can learn the short-term and long-term feature representation from the mel-scale filter banks [169]. In addition, to balance the positive and negative samples, a random sampling approach is adopted in the model training stage before using LSTM. Using the DepAudioNet, different scale representations, i.e., high-level, short-term, and long-term features, are extracted. To further explain different representations between healthy controls and depressed subjects, Fig. 6

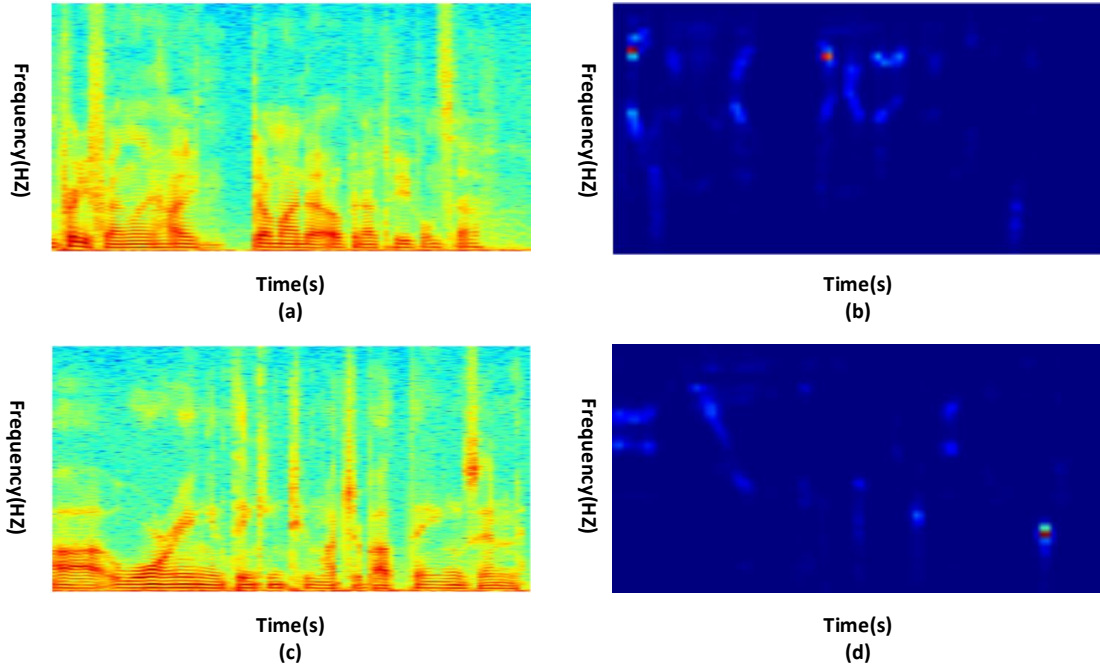


Figure 6: The visualization of spectrogram and mel-scale filter banks. (a) and (b) represent the spectrogram and filter bank features from an audio segment of a health control. (c) and (d) show the spectrogram and filter bank features from an audio segment of a depressed individual [23].

provides a comparison of the spectrogram and filter bank features extracted from an audio segment. The goal of the authors is to try to use deep learning methods to estimate the severity of depression. Most importantly, despite the small size of the training data, the deep learning methods can also learn the discriminative patterns from audio signals.

Even though the used depression databases have only a limited number of samples, deep learning-based depression recognition approaches aroused great attention from numerous researchers. Then in 2018, a fusion of deep-learned and hand-crafted features was used, capable of effectively measuring the severity of depression from speech. In the framework, DCNN was used to learn and fuse the shallow and deep patterns for evaluating the severity of depression. Specifically, LLD features are extracted by OpenSMILE toolkit [170] as hand-crafted features from audio. Median robust extended local binary patterns (MRELBP) are extracted as hand-crafted features from spectrograms. Raw audio and spectrograms are used as input into the DCNN to obtain the deep learned features. To learn the complementary representations between the hand-crafted features and the deep-learned features from the Raw-DCNN and Spectrogram-DCNN, a joint fine-tuning technology is used. In addition, to overcome the issue of the small number of samples, a data augmentation approach is introduced. Most importantly, the proposed scheme presents an end-to-end architecture for depression recognition [33] (see Fig. 7). The contribution of the work [33] is the attempt to fuse hand-crafted and deep learned features from speech for depression estimation. Also, the authors [33] extracted texture features from spectrograms for predicting the severity of depression. The authors validated the proposed method and obtained a good performance on AVEC2013 and AVEC2014 databases with RSME of 10.00 and 9.98, respectively (see Table 4).

Due to the limited sizes of databases available for depression recognition, different studies have proposed augmenting the data somehow. For instance, in [171], Deep Convolutional Generative Adversarial Network (DCGAN) was proposed to augment the size of data samples to enhance the accuracy of ADE task from audio signals. To validate the performance of the augmented features, three measurement criteria have been proposed, i.e., spatial, frequency, and representation learning. The proposed architecture was capable of achieving comparable performance with most of the methods on the DAIC database, with RMSE of 5.52 and MAE of 4.63 (see Table 4). As illustrated in Fig. 8, the DCGAN framework contains a learning strategy with two levels, to improve the convergence speed of the training. In the first level, feature maps were split into 9 blocks with a size of 28×28 . For every block, a DCGAN model is used

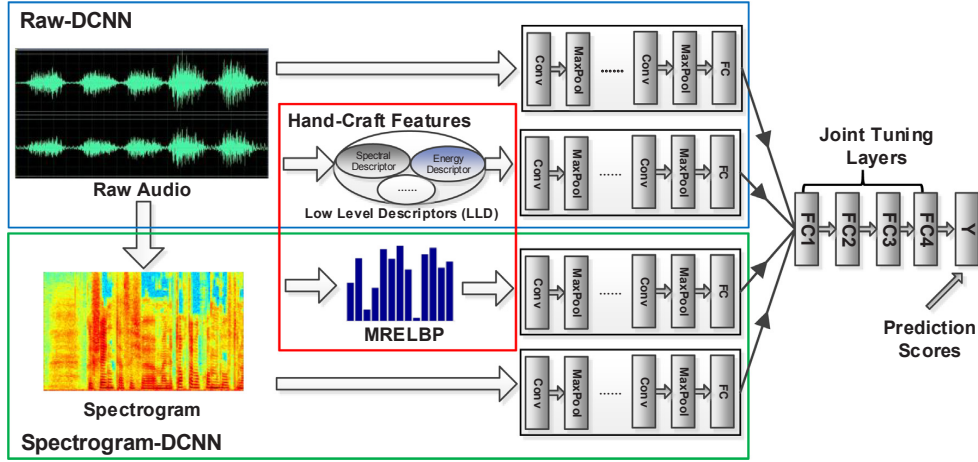


Figure 7: The framework for depression proposed in [33]. The four-stream deep features, i.e., hand-crafted features (LLD, MRELBP), deep-learned features (raw audio, spectrogram), are fused for deep depression recognition. Raw-DCNN (Top) adopts LLD and raw audio signals as input, whilst Spectrogram-DCNN (Bottom) utilizes MRELBP and spectrogram features as input. The red box represents Hand-Crafted features. The other two arrows denote deep-learned features. The BDI-II score is calculated by averaging and aggregating the output of the four DCNN branches.

to represent synthetic representations. After that, 9 DCGANs are generated with the same architecture. The output of the first level (with the size of $9 \times 28 \times$) is fed into the second level to obtain the global features. The advantage of this architecture is that complex training is transformed to a more straightforward procedure.

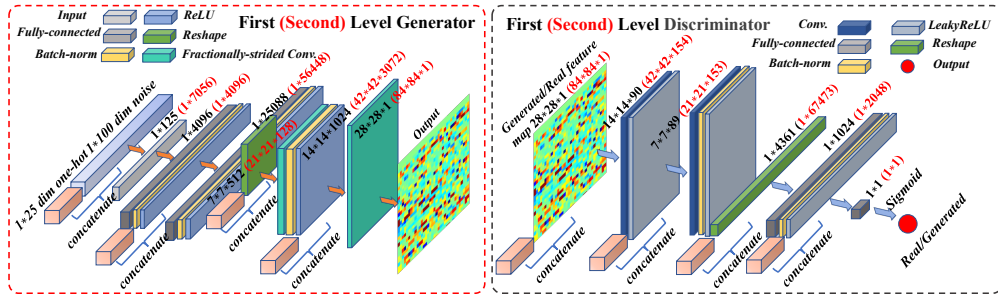


Figure 8: An illustration of the framework proposed in [171]. In the DCGAN framework, a 2-level hierarchical learning strategy is proposed to streamline the training procedure [171]. In the first level, feature maps are split into 9 blocks with size 28×28 . For every block, the model generates synthetic representations. After that, 9 DCGANs are generated with the same architecture. The output of the first level (with the size of $9 \times 28 \times 28$) is fed into the second level to obtain the global features.

Table 4 lists the reviewed approaches for depression recognition with audio cues. Table 4 shows that many studies adopted the DCNN model to extract deep features for predicting the depression scale. From the algorithmic perspective, researchers adopted the widely used DL technologies for evaluating the severity of depression, i.e., DCNN, LSTM, etc. It is noteworthy that the raw audio signals are fed into the DCNN directly to overcome the drawbacks of the conventional feature design methods [33]. Furthermore, Niu et al. [163] attempted to convert the audio segments into a spectrogram to feed into the deep architecture. They sampled the audio clips at 8KHZ, and adopted STFT using a Hamming window with 32ms and shift with 16ms to generate 129-dimensional spectrograms on the two AVEC2013 and AVEC2014 databases. In addition, they found that the optimal spectrogram length is 64 frames (1s) with the shift of 32 frames (0.5s) of the two databases, respectively.

In 2021, Niu et al. [172] proposed a novel framework that integrates with the Squeeze-and-Excitation (SE) component and a Time-Frequency Channel Attention (TFCA) block to represents informative characteristic related to depression. Additionally, to consider the time-frequency features of the data, a Time-Frequency Channel Vectorization (TFCV) block is proposed to form the tensor. Moreover, they integrated the introduced blocks (i.e., TFCA and

TFCV blocks) and the two blocks (i.e., Dense block and Transition Layer) of the DenseNet into a unified framework to generate Time-Frequency Channel Attention and Vectorization (TFCAV) network. The contribution of this work [172] is that time-frequency attributes are considered to learn the informative patterns from spectrograms. The performances of the introduced method obtained on AVEC2013 and AVEC2014 with RMSE of 8.32 and 9.25, respectively.

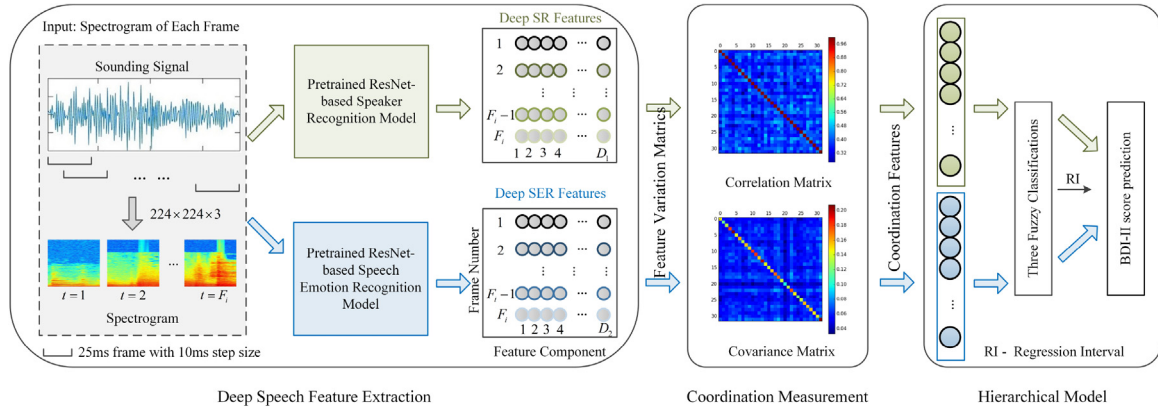


Figure 9: Illustration of the architecture for ADE from speech proposed in [173]. The architecture can be divided into three steps: deep speech feature extraction, coordination measurement, and hierarchical model construction. In the first step, a spectrogram is used for learning the features of the frame-level speaker recognition (SR) and speaker emotion recognition (SER) from the pre-trained SR and SER models. In the second step, the FVCM algorithm is adopted to compute the correlation and covariance coefficients of the time-delayed multi-channel variations to obtain the coordination features. In the third step, a ADE model is introduced.

In [173], the authors proposed a deep architecture for ADE from speech with two contributions. The first one is that Speaker Recognition (SR) and Speaker Emotion Recognition (SER) features were fused to improve ADE’s performance. The second contribution is that the Feature Variation Coordination Measurement (FVCM) algorithm is used to model the correlation and covariance coefficients of the time-delayed multi-channel variations (see Fig. 9).

Fig. 10, (a) and (b) illustrate the audio and video sub-sequences for a healthy control. (c) and (d) represents the same patterns for a depressed individual. Imagesc toolkit of MATLAB is adopted to plot the figures. From the Fig. 10, one can see that there are difference between (a) and (c) or between (b) and (d). For instance, the blocks enclosed by red rectangles are discriminative between the individuals, whereas the parts covered by green rectangles more are similar. Based on Fig. 10, we can make the following observations: The discriminative patterns of audio and video frames have different contributions for healthy controls and depressed subjects. The mentioned studies provide motivation for the subsequent works in ADE.

4.3. Deep ADE Networks for Video Modality

Besides the audio modality, visual cues also important for deep depression recognition. Therefore, various researchers from the affective computing field have explored the discriminative patterns in videos for ADE. In the following, we introduce the studies based on video for ADE. Accordingly, we divide the mentioned methods into two groups based on the input of the deep networks: deep ADE networks for single image and deep ADE networks for image sequences.

4.3.1. Deep ADE Networks for Single Image

The study [28] was an initial attempt to adopt deep learning for depression detection from static images. A two-stream network was developed in their proposed framework to use facial images and optical flow features to learn the depression patterns (Fig. 11). Appearance-DCNN and Dynamics-DCNN have been introduced to model the static and dynamic patterns for depression recognition. The Appearance-DCNN includes two steps. The first step consists of training a model from scratch on a public CASIA WebFace Database with 494,414 images from 10,575 subjects [174]. After that, the deep model includes the discriminative representations related to facial structures, which can provide sufficient information for the ADE task. However, the pre-trained model can not be directly used for ADE. The second step is to fine-tune the pre-trained model for ADE. However, the ADE task based on AVEC2013 and AVEC2014 can

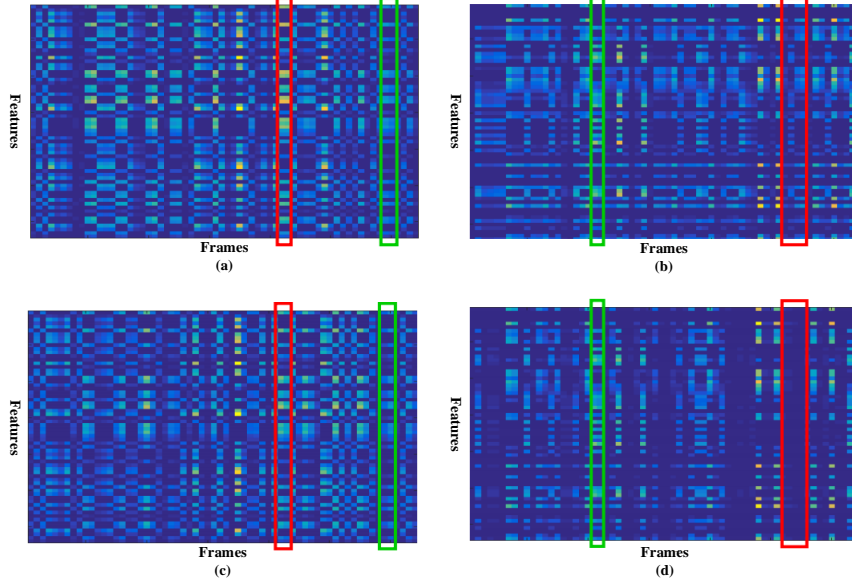


Figure 10: The example of audio and video features for health controls and depressed subjects. (a) and (b) represent the audio and video characteristic for the healthy controls (No. 203-1) of the AVEC2013 database, respectively. (c) and (d) represent the patterns for the depressed subjects (No. 236-1) of the AVEC2013 database. The scale of depression of No. 203-1 is none (3), while No. 236-1 is 23 (moderate). To further illustrate the depression scales, red and green rectangles are used to represent discriminative and less discriminative feature vectors.

be viewed a regression problem from the perspective of machine learning. Hence, the softmax loss function is changed into the Euclidean loss for ADE. To further model the dynamics between several consecutive video frames, optical flow displacements are computed for the Dynamics-DCNN. The subtle dynamic patterns and motions of the face are explored, and redundant information of videos is reduced by using the optical flow. In particular, the study leverages the ability of the existing large models to predict the BDI-II scores on small datasets. Most importantly, this work [28] has given a certain inspiration for the following works based on deep learning for depression recognition and analysis. The detailed architecture of Appearance-DCNN and Dynamics-DCNN of Fig. 11 for ADE.

In 2018, Zhou et al.[175] proposed a novel deep architecture named *DepressNet* to learn representations from images for depression recognition, as shown in Fig. 13. Different deep architectures (AlexNet, ResNet, GoogleNet) were pre-trained on the CASIA database. *DepressNet* is constructed by changing the softmax layer into a regression layer, followed by a global average pooling (GAP) layer, as shown in Fig. 13. Specifically, the *DepressNet* consists of four bottleneck blocks, which include 3, 4, 6, and 3 bottleneck structures. After that, the deep model is fine-tuned on the AVEC2013 and AVEC2014 databases. Meanwhile, the loss function also changed into the squared loss, which can be written as follows:

$$L = \frac{1}{2M} \sum_{j=1}^M (g(y_j) - \ell_j)^2, \quad (1)$$

where M represents the batch size, $g(y_j)$ and ℓ_j are the predicted value and the label of the j -th face image of sample x_j , respectively. Another method, a multi-region *DepressNet* (MR-*DepressNet*) [175] has been designed for learning different scale models for overall depression recognition, as shown in Fig. 14. In this architecture, to learn the discriminative patterns from different regions and full images, a four-stream *DepressNet* is proposed. In order to learn more robust representations, the output of the four sub-architectures is combined at the cost function layer. Formally, the loss function of MR-*DepressNet* can be written as:

$$L = \frac{1}{2M} \sum_{j=1}^M \left(\frac{1}{J} \sum_{j=1}^J g_j(y_j) - \ell_j \right)^2 \quad (2)$$

Where M represents the batch size, J represent the number of the image regions, $g_j(y_j)$ denotes the output of the

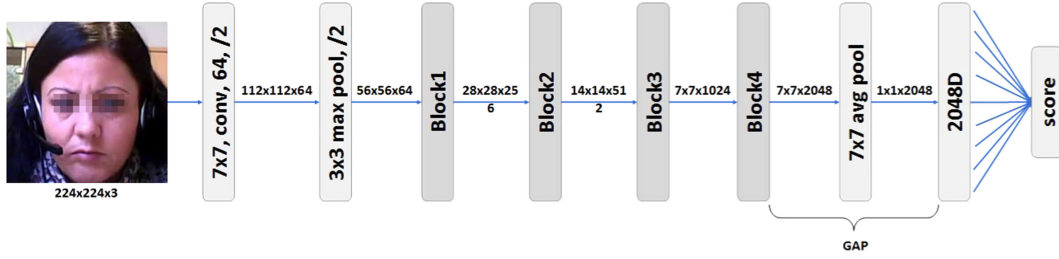


Figure 13: The detailed architecture of DepressNet for ADE task [175]. In this architecture, the facial images are first pre-processed by the OpenFace toolkit to ensure they have the same scale. The architecture has residual connections, similar to those in the popular ResNet architecture. DepressNet contains four blocks, which consist of 3, 4, 6, 3 bottleneck architectures for feature representation. Then the 2048D features are extracted from the architecture for ensemble depression prediction.

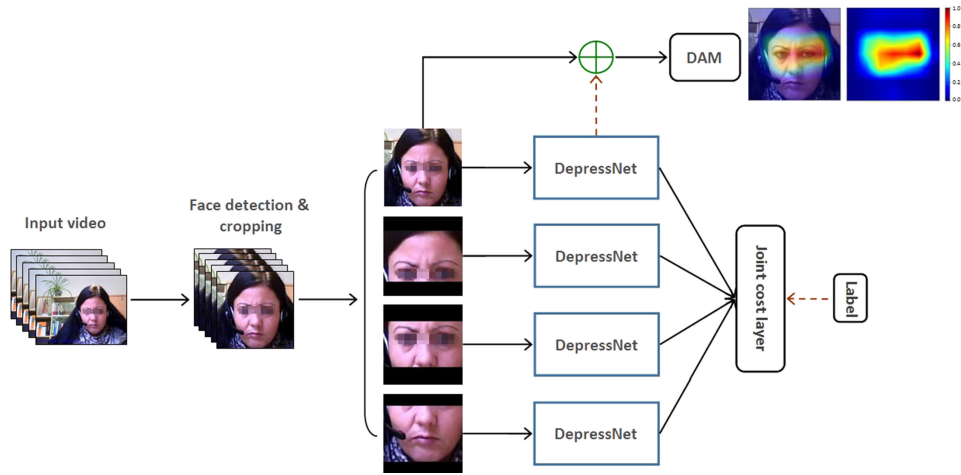


Figure 14: The detailed architecture of Multi-Region DepressNet for ADE task from [175]. In this architecture, the facial images are first pre-processed by OpenFace toolkit to ensure them at the same scale. Then the facial area is divided into different regions, which are fed together with full face into DepressNet to estimate the BDI-II score.

1) They leverage the large scale database (e.g., CASIA, VGG, etc.) to pre-train the deep models by using the deep architectures (e.g., GoogleNet, VGG, ResNet, etc.). 2) They improve the performance of the deep models by fine-tuning them on depression databases, e.g., AVEC2013 and AVEC2014, etc. 3) Moreover, some works try to improve the performance of depression recognition via designing a particular loss function for the ADE task.

Interestingly, He et al. [177] introduced a novel network that combined a 2D-CNN networks and attention mechanism for depression recognition. In this method, the authors proposed an integrated architecture – Deep Local-Global Attention Convolutional Neural Network (DLGA-CNN) for ADE, which utilizes DCNN with attention mechanism, and weighted spatial pyramid pooling (WSPP) to model a global feature. Two branches are designed: Local Attention-based CNN (LA-CNN) concentrates on the local patches, while Global Attention-based CNN (GA-CNN) models the global features from the entire facial region. In order to learn the complementary patterns from the two branches, Local-Global Attention-based CNN (LGA-CNN) is introduced. After the aggregation of features, WSPP is adopted to extract the depression representations. More importantly, compared to the previous methods, the proposed method did not leverage a large-scale database to pre-train the deep model, but rather is considered as an end-to-end scheme for ADE (see Fig. 18).

Table 4 suggests that most studies using DL until now have adopted the DCNN architecture to predict the depression scale. In addition, attention mechanism [191], has also been utilized in depression recognition [177]. As for pre-processing, the researchers have mainly used MTCNN, OpenFace, Dlib toolkits to detect and crop the facial region to lay a solid foundation for depression detection.

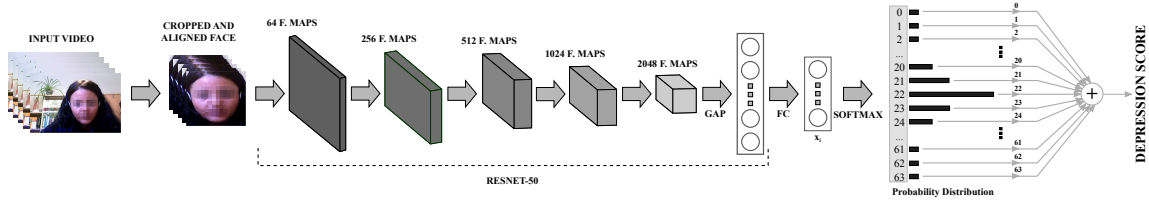


Figure 15: The proposed method [30] is to estimate the severity of depression. The videos were first processed to generate the aligned facial images. Then ResNet-50 was fine-tuned to extract the discriminative features and follows a GAP layer to pool the features. Finally, expectation loss is used to weight the performance of the proposed method.

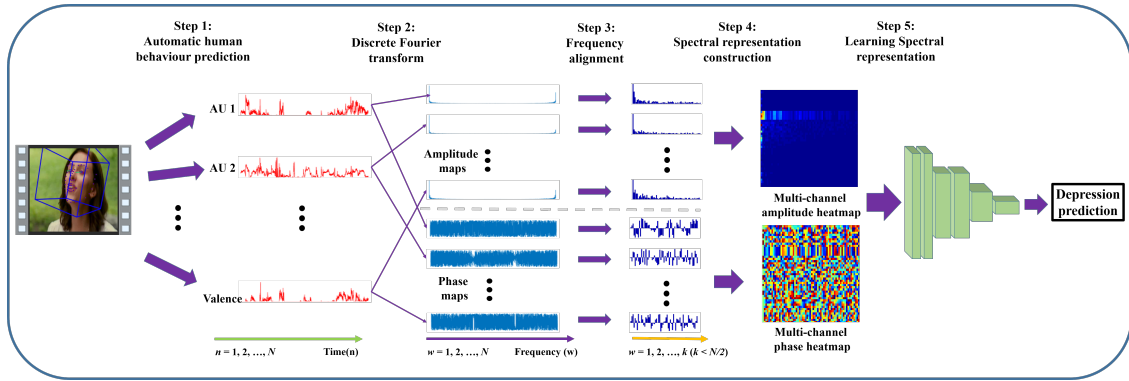


Figure 16: The proposed method [31] to estimate the severity of depression, which can be divided into five steps: 1) Multi-channel human behaviour primitives were extracted from videos; 2) Human behaviour primitives were then transformed into spectral signals with multiple frequency patterns from all frames; 3) Owing to the symmetric of spectral signals, high-frequency pattern were removed to retain the discriminative information of human behaviours from the videos; 4) Multi-channel amplitude heatmap and multi-channel phase heatmap are constructed from spectral signals; 5) DCNN and Artificial Neural Networks (ANNs) method were used to predict the depression scales.

4.3.2. Deep ADE Networks for Image Sequences

Though discriminative patterns based on single image features have been widely adopted in ADE tasks, and have obtained promising performance, these works nonetheless neglect the temporal information that could be useful for the ADE task. To resolve this Jazery et al. [27] proposed using C3D and RNN to extract spatiotemporal features in two different scales from the video clips for depression recognition. The proposed framework consists of the two components, i.e., loose and tight scale feature extraction components, which use fine-tuning of deep models and temporal feature aggregation. The C3D Tight-Face model is utilized for learning of a tight (i.e. high-resolution) features, while the C3D Loose-Face model is trained on larger face regions to learn the global features. Then an RNN is adopted to model the temporal features learned by the C3D Tight-Face and C3D Loose-Face models. Finally, a mean operation is utilized for prediction. The main contribution of this work [27] is the temporal framework that learn facial features on different scales. Furthermore, different feature aggregation stages can combine the features from different scales, which can benefit depression scale prediction (see Fig. 19).

Subsequently, Melo et al. [176] proposed a combination of different C3D architectures to learn Spatio-temporal patterns from full-face and local regions, and further combined them with 3D Global Average Pooling (3D-GAP) for predicting depression. The local C3D architecture learns discriminative information of the eye region, while the global C3D architecture focuses on learning the spatio-temporal patterns based on the whole facial region. In addition, 3D-GAP is also used to aggregate the spatio-temporal features from the last convolutional layer (see Fig. 20). The proposed method was assessed on the AVEC2013 and AVEC2014 databases, and it obtained an improved performance with the RMSEs of 8.26 and 8.31 as compared with the state-of-the-art methods, respectively (see Table 4).

Uddin et al. [26] used LSTM to model the sequence information from video data. Moreover, deep facial expression features were extracted by a deep CNN and then pooled by a Temporal Median Pooling (TMP) technology to feed the LSTM module for ADE. Various experiments were performed on the two datasets (AVEC2013 and AVEC2014),

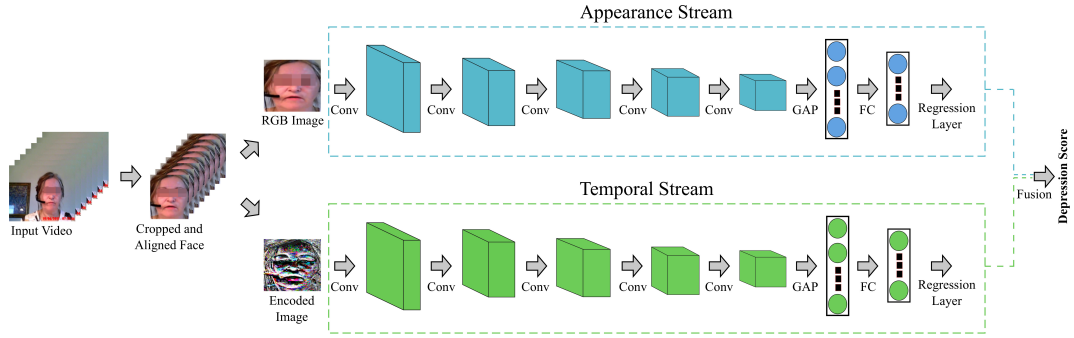


Figure 17: This method [176] is proposed to estimate the severity of depression. The appearance stream takes the static images as input, while the temporal stream takes image sequences as input. A simple fusion method, i.e., average pooling, was used to fuse the outputs of the two networks for the ADE task.

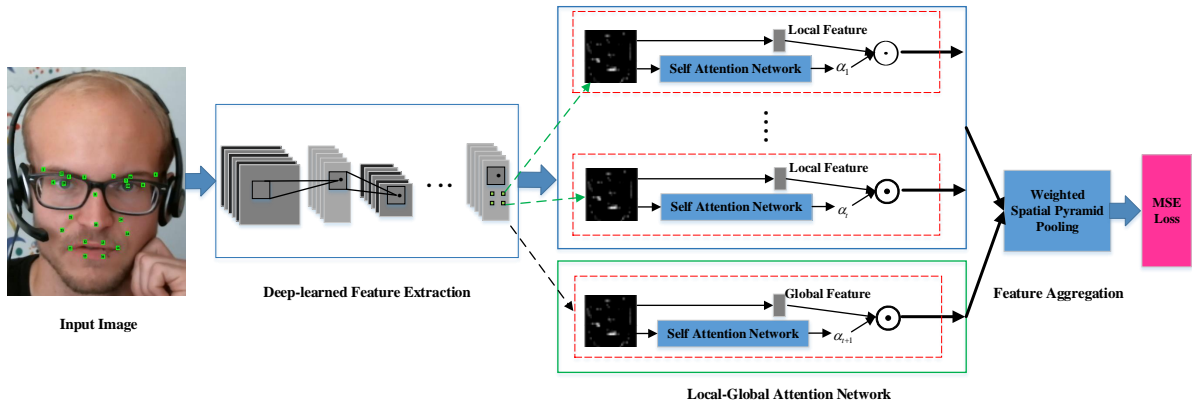


Figure 18: The proposed method DLGA-CNN for ADE [177]. The facial image is obtained by OpenFace toolkit [165]. Then a typical DCNN is designed for feature representation to generate a discriminative the feature maps. To extract informative features, the local and global self-attention networks are designed. To obtain scale-invariant feature representations over multi-scale feature maps, WSPP is used. In addition, two fully connected layers and a Mean Square Error (MSE) loss layer are adopted for ADE.

indicating that the proposed methodology surpasses most existing methods (see Table 4). The contribution of work [26] is that a Volume Local Directional Number (VLDN) dynamic feature is designed to model the trivial emotions from facial regions. In [181], a novel 3D framework, the multi-scale spatiotemporal network (MSN), was developed to learn the characteristic information of the video clips. In the work, several parallel convolutional layers were used to learn considerable spatio-temporal variations from facial expressions. The model adopts several receptive fields to maximize the exploitation of distinct spatial areas from facial region for ADE (see Fig. 21).

In 2021, several works [182, 183] have proposed to predict the severity of the depression. In [182], the authors proposed an end-to-end intelligent system to generate discriminative representations from the entire video clip. Specifically, a 3D-CNN combined with a Spatiotemporal Feature Aggregation Module (STFAM) is trained from scratch on AVEC2013 and AVEC2014 data, which can learn the informative patterns of depression. In the STFAM, channel and spatial attention mechanism as well as an aggregation method, namely 3D DEP-NetVLAD, are integrated to capture the compact characteristic based on the feature maps. Case studies are introduced to describe the applicability of the proposed intelligent system for ADE (see Fig. 22).

In [183], a new DL architecture named Maximization and Differentiation Network (MDN) is proposed to model the variations of facial expressions closely related to depression. The MDN is designed without 3D convolutions, and it exploits discriminative temporal patterns learned by two different blocks that model either smooth or sudden facial variations. Finally, they designed the models with 100 and 152 layers and validated the deep models on the

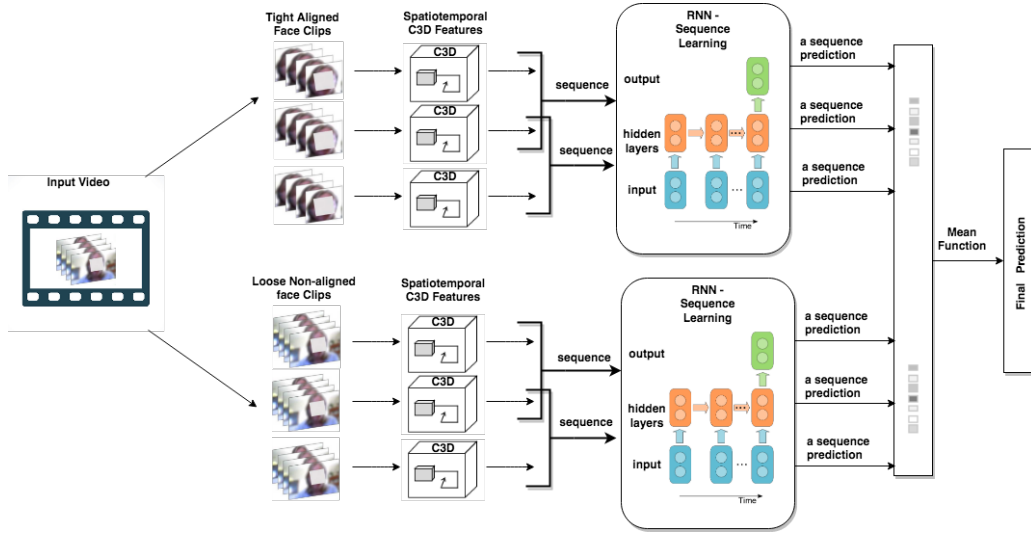


Figure 19: The pipeline of the framework for estimating the scales of depression with Deep C3D and RNN from videos proposed in [27]. Discriminative features are extracted at two different scales. C3D Tight-Face model learns a tight (i.e. high-resolution) feature representation, while the C3D Loose-Face model is trained on larger face regions to learn global features. An RNN is adopted to model the temporal features based on outputs of the C3D Tight-Face and C3D Loose-Face models. Finally, a mean operation was used to generate the predictions.

AVEC2013 and AVEC2014 databases (see Fig. 23, Fig. 24). The proposed model obtains competitive RMSEs of 7.55 and 7.65 on the AVEC2013 and AVEC2014 databases, respectively (see Table 4).

Based on the works mentioned in this section (see Table 4), we can make the conclusions as follows:

1. In comparison with static features, image sequences can capture short-term and long-term Spatio-temporal information from videos. This can improve training of a deep discriminative model for ADE task.
2. From the perspective of training, most of the works include pre-training and fine-tuning stages for ADE. To date, there does not exist an end-to-end scheme for ADE from image sequences.
3. Summarizing the results in this section (Section 4.3.2), most of the works obtained a comparable performance. So far, the method published in [181] obtains the best result with RSME of 7.55 on AVEC2013, and the method from [183] gets the best result with RMSE of 7.61 on AVEC2014, respectively.

4.4. Deep ADE Networks for Multi-modal Fusion

Apart from the above-mentioned single modalities audio (see Section 4.2) and video (see Section 4.3), multi-modal fusion methods can enhance the performance of depression prediction. A combination of DCNN and DNN methods was proposed in 2017 for ADE [118], using different models to combine audiovisual features and textual inputs from a transcript. For each respective single modality, hand-crafted features were fed into a DCNN to model global-scale features, then input into a DNN to assess the PHQ-8 scores. To promote the performance of depression recognition, a multi-modal fusion scheme was formulated. Subsequently, the three single models (audio, visual, text) were fused together and input into a DNN to predict the severity of depression defined by the PHQ-8 depression scale. Moreover, Paragraph Vector (PV) was proposed to learn the distributed representations for the text descriptors. Besides, a novel video feature was proposed, i.e., Histogram of Displacement Range (HDR), capable of learning the displacements and speed of facial landmarks. Experiments were performed on the AVEC2017 challenge. It obtained comparable performance, with the RMSE of 5.97 and the MAE of 5.16 on the test set (see Table 4). In [118, 186], a hybrid depression recognition framework based on audiovisual and text descriptors was proposed. In the framework, DCNN and DNN were first used to classify depressed subjects and healthy controls. In the studies [184, 119], the method mentioned in [186] was also adopted to predict the severity of depression, with promising performance (see Fig. 25).

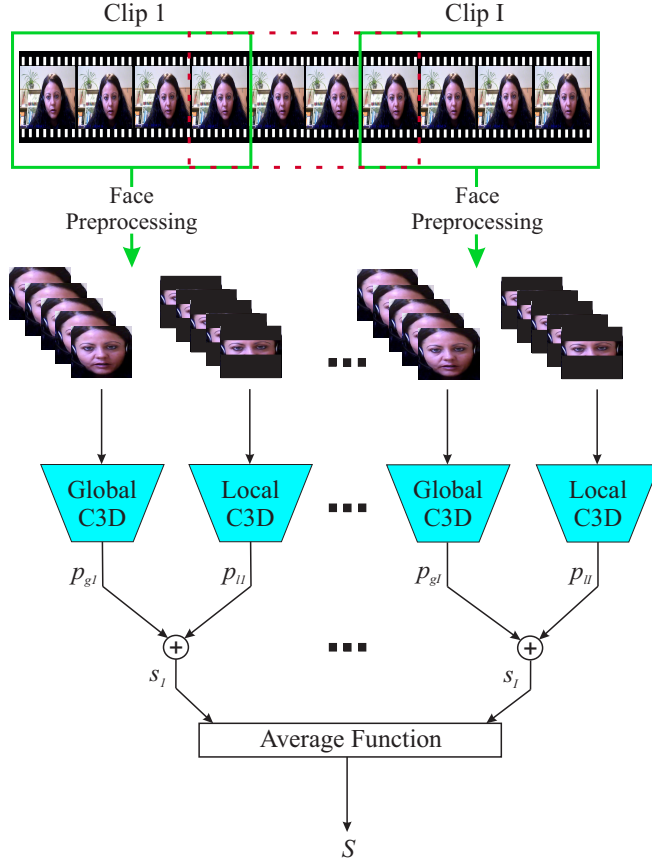


Figure 20: The pipeline of the proposed scheme for ADE with local and global C3D from videos [176]. The video clips were pro-processed by the MTCNN toolkit. Then the two discriminative features at different scales based on C3D were extracted and concatenated. After that, the average function was used to pool the final features for predicting the depression scores.

Meanwhile, a novel Bipolar Disorder Corpus was released for academic research [145] and then used for the AVEC2018 Bipolar Disorder Sub-Challenge. Based on the AVEC2018 database, in [193], a novel architecture fusing a DNN and Random Forest was proposed for bipolar depression analysis. In [194], to address bipolar disorder (BD) with irregular variations among different episodes, a new architecture, i.e., IncepLSTM, was designed, capable of combining Inception module and LSTM of feature sequence with learning multi-scale temporal patterns for BD analysis. Experiments were performed on the AVEC2018 dataset, demonstrating the efficiency of the proposed method. In addition, other works also adopted conventional machine learning methods for BD recognition [195, 196]. However, the AVEC2018 database has not been extensively employed by researchers from the affective computing community so far. Notably, [185] presented a novel method integrating unsupervised learning, transfer learning, and hierarchical attention from speech to assess the depression scale (see Fig. 26). The proposed method was evaluated on the AVEC2017 depression challenge, and RMSE and MAE are 5.51 and 4.20, respectively (see Table 4).

In addition, two feature sets were introduced to determine the duration of sequential landmarks. As reported by extensive experiments evaluated on DAIC-WOZ and SH2 databases, the model can learn the patterns effectively as compared with the previous speech-based depression recognition methods [197].

To learn the auxiliary information between audio and video cues, a new Spatio-Temporal Attention (STA) architecture and a Multi-modal Attention Feature Fusion (MAFF) method was proposed to extract the multi-modal features from audiovisual cues for predicting the depression scale, i.e., BDI-II score. The proposed method comprises 2D-CNN, 3D-CNN, and an attention mechanism to learn the deep features for depression scale prediction. Extensive experiments are carried on the AVEC2013 and AVEC2014 databases, demonstrating that the proposed deep architec-

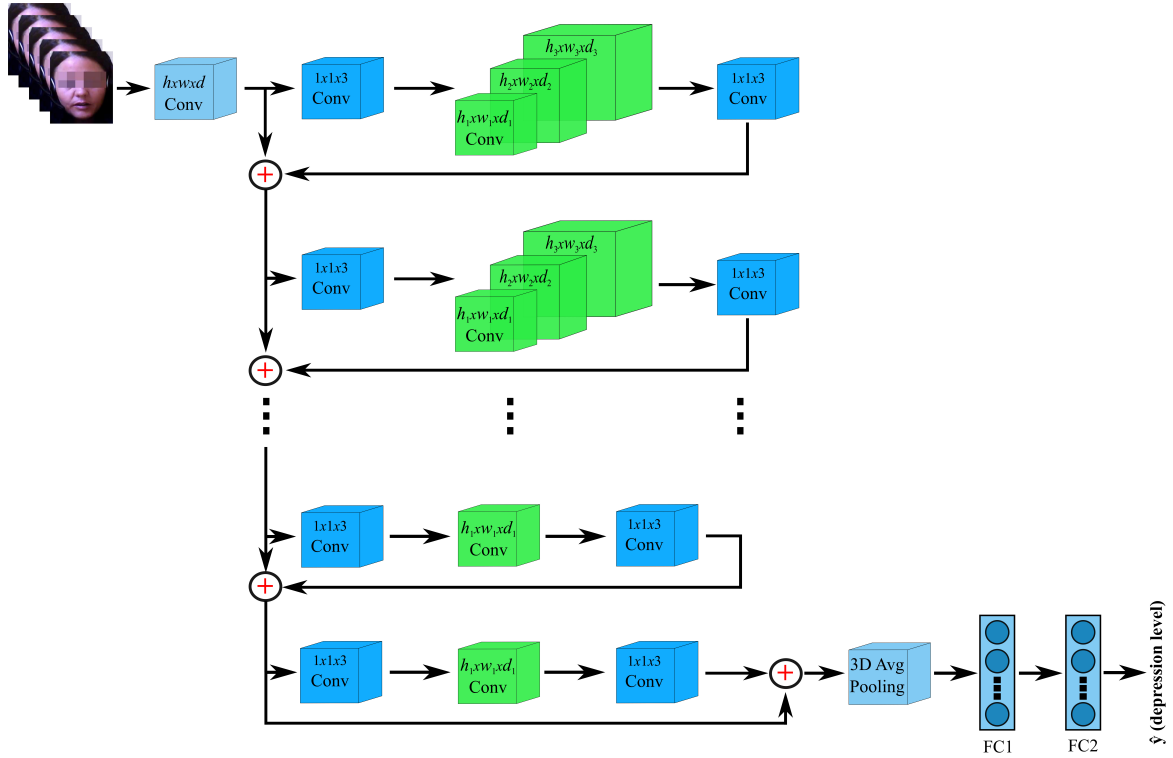


Figure 21: The illustration of the proposed MSN for ADE [181]. In the framework, several parallel convolutional layers are used to learn considerable spatio-temporal variations from facial expressions. The model adopts several receptive fields to capture the multi-scale pattern of depression for ADE.

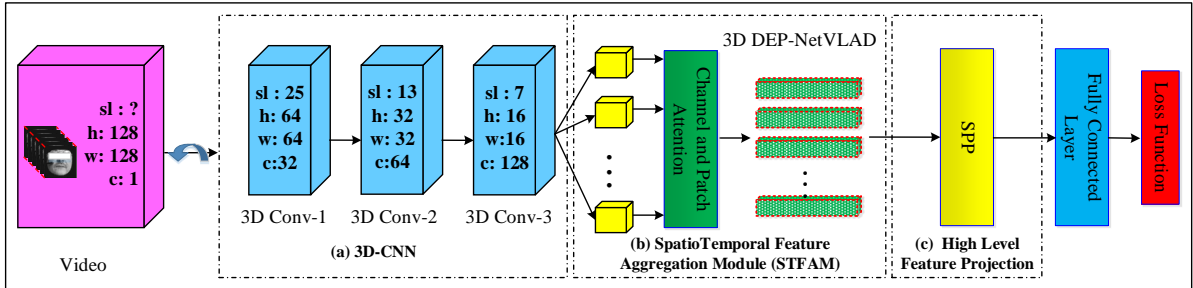


Figure 22: The illustration of the proposed architecture for ADE [182]. In the framework, which consists of the following steps. The first step is that facial images are cropped and aligned by the OpenFace toolkit. The second stage is that 3D-CNN is adopted to extract the local and spatial-temporal feature representations related to the symptoms of depression. The third stage is that Spatio-temporal Feature Aggregation Module (STFAM) is used to aggregate the discriminative features over the local features. An SPP layer is adopted to represent multi-scales on top of the output of STFAM. Finally, a fully connected layer and MSE loss function is used to predict the final BDI-II score.

ture outperformed most of the existing studies [163] (see Fig. 27). To sum up, the mentioned works appeared for the image sequences for depression recognition, and they leverage mature DL technology (e.g., DCNN, RNN, LSTM) to learn the deep discriminative patterns for depression estimation. In addition, an attention mechanism is also used to learn the salient patterns from the deep-learned features. Moreover, in the Detecting Depression with AI Sub-Challenge (DDS) of AVEC2019, several works are also focused on adopting AI technology to estimate a subject's depression scale (Table 4).

Also, the authors in [180] used the LSTM to model the interactions in audio and text features based on the

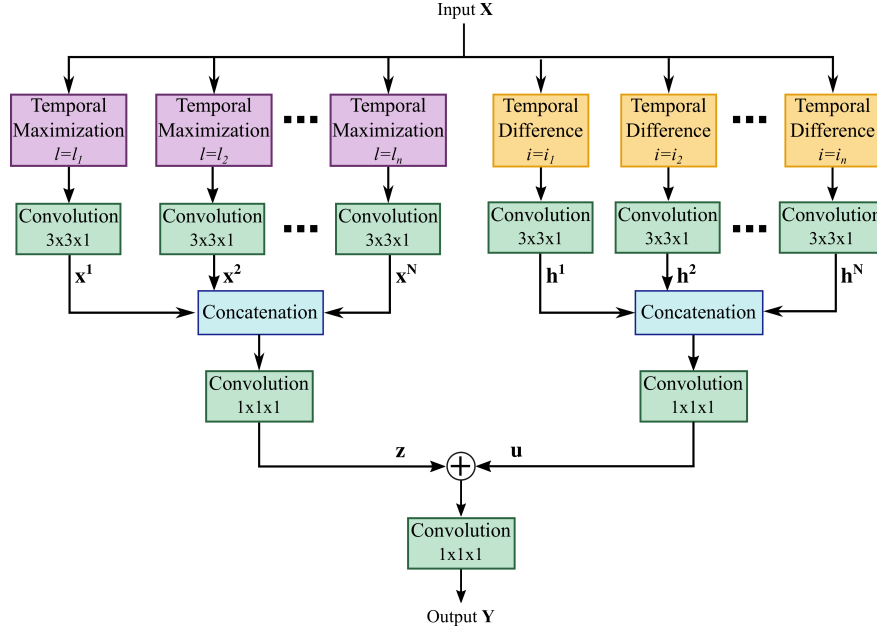


Figure 23: A detailed illustration of the MDN module [183]. The module consists of two blocks. The left block is a maximization block to model spatio-temporal patterns. $\mathbf{X} \in \mathbb{R}^{N \times T \times H \times W \times C}$ is the feature map, in where N, T, H, W and C represents the the batch size, temporal depth, height, width, and the number of channels, respectively. l_1, l_2, \dots, l_N represent the branches of the maximization block. x_i represents the output of i -th branch. $z = \mathcal{H}\{\bigcup_{n=1}^N x^n\}$ is the output of the maximization block. The temporal difference block (right) learns spatio-temporal variations. i_1, i_2, \dots, i_N represent the branches of difference block. h_i represents the output of i -th branch. $u = \mathcal{H}\{\bigcup_{n=1}^N h^n\}$ is the output of the difference block. Finally, the two blocks are combined to obtain the final output.

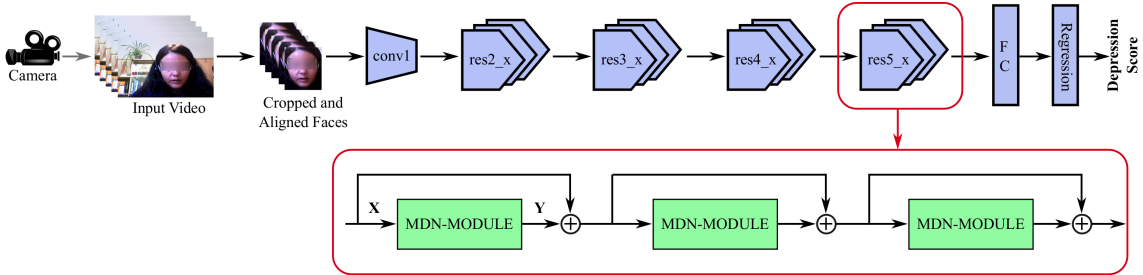


Figure 24: Illustration of the MDN architecture [183]. Firstly, facial images are cropped and aligned by MTCNN. Secondly, 3D residual networks [192] are pre-trained on VGGFace2 dataset for image classification. Thirdly, the 3D residual deep model combined with the MDN module is fine-tuned on AVEC2013 and AVEC2014 databases to compute the BDI-II scores.

AVEC2017.

To summarize our findings (see also Table 4):

- From the perspective of modality, the multi-modal fusion method yielded the optimal performance for ADE on every database. On the AVEC2013 and AVEC2014 databases, Niu et al. [163] achieved the best accuracy with RMSE of 7.03, and MAE of 5.21. On the DAIC database, the method proposed in [118] obtained the best performance with RMSE of 5.40, and MAE of 4.35, respectively. Though multi-modal fusion yields the optimal performance for ADE, this method is very complicated when fusing the complementary information between audio and video cues. Consequently, a wide range of researchers have focused on the video modality to learn the discriminative patterns around facial regions, which is probably explained by the success of computer

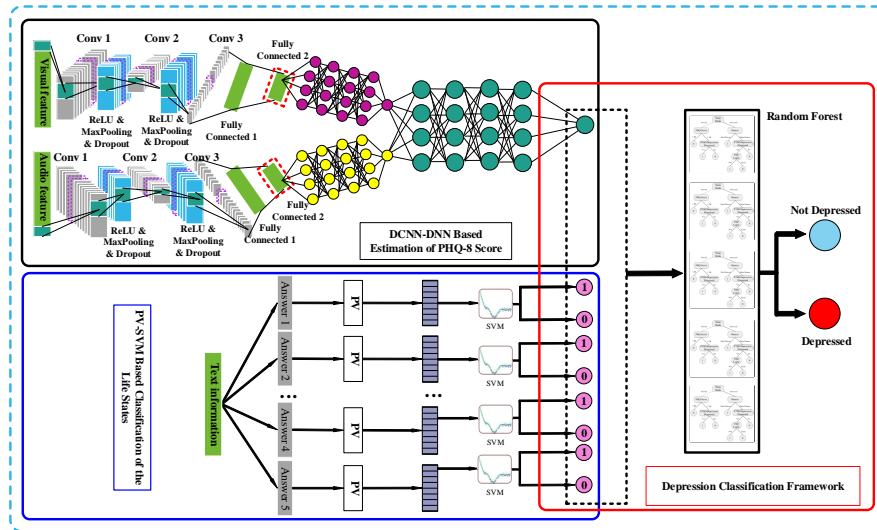


Figure 25: The integration of shallow and deep architecture for multi-modal depression prediction and classification [119]. The method includes three modules: 1) The audiovisual DCNN-DNN prediction module. DCNNs learn high-level features from hand-crafted audio/video features from audiovisual signals. Then PHQ-8 scores were predicted by inputting the high-level features to a DNN, and then fed to a DNN to get the final prediction (surrounded by the black rectangle). 2) The classification module. PV-SVM model is adopted to check the presence or absence of the psychoanalytic symptoms, e.g., sleeping disorder and feelings (the blue rectangle). 3) The depression classification module. To get the final classification results based on the results of 1) and 2), a random forest method was adopted to classify the participants as the healthy controls and depressed subjects (the red rectangle).

vision and DL in general.

- From the database point of view, AVEC2013, and AVEC2014 have obtained the most attention. The reason is that the audio and video clips are contained in the AVEC2013 and AVEC2014 databases. Thus, the researchers can leverage DL technology to learn the compact representation from the video clips. For the DAIC database, the database organizers only provide the audio data samples, limiting their use for ADE.
- From the perspective of DL, the DCNN is commonly used by different works to learn the discriminative patterns from both static images and hand-crafted features. To model the sequential information from the video sequences, a 3D-CNN is also used in many works. Different variations based on 3D-CNN are also proposed to predict the severity of depression [183].

Furthermore, since 2015, there also exist works that have not adopted DL technology for depression estimation, (e.g., [198, 33, 199, 200, 201]). In particular, Sadari et al. [200] used ordinal logistic regression for depression recognition and proposed a new way for the task. In addition, based on the database of AVEC2017, numerous methods have been proposed for depression recognition. In [202], the authors analyze the relationship between the ground truth and predictions when measuring the severity of depression. They design a system validated on the AVEC2017 depression database. They found that depression recognition is an ordinal problem. Also, He et al. [201] introduced a promising feature descriptor, i.e., median robust LBP-TOP (MRLBP-TOP), that can learn patterns on different scales from image sequences. Dirichlet process FV (DPFV) has also been proposed to learn the global patterns from the segment-level features. In addition, Bipolar Depression (BD) has also aroused attention in the affective computing field. According to Table 5, notably, various works have been considered for estimating the BD. Researchers still continue using DCNN, LSTM and DNN to extract the deep features to represent the severity of BD from the perspective of methods.

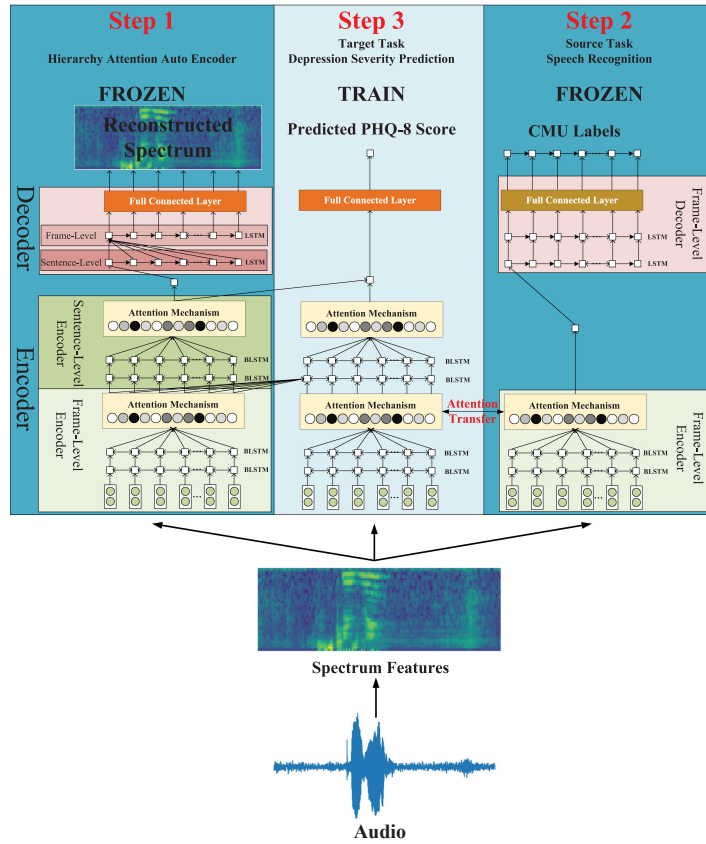


Figure 26: A combination of hierarchical attention and auto-encoder architecture was proposed for depression recognition [185]. First, a hierarchical attention auto-encoder network was trained based on a bottom-up scheme, in which an encoder with the attention mechanism generates a sentence-level encoder to construct the hidden vectors of the entire context. Then the vectors are performed at sentence-level attention to generate a latent representation of the clinical interview. Second, the representation was input into a decoder to re-capture the input feature. Third, the representations of the sentence are fed into Bi-LSTM to assess the PHQ-8 scores. After that, the parameters of the auto-encoders were frozen. Then the obtained model captures its attentions by a speech recognition task, and transformed into the hierarchical depression detection system for ADE.

5. Open Issues and Promising Directions

This section introduces open challenges in ADE and suggests promising future directions. We want to encourage the affective computing community to study together on the mentioned issues to boost ADE development. Consequently, our goals are: (1) to promote AI-based ADE frameworks to apply in real life, especially in hospitals, psychiatric centers, and so on. This comprises assessing and polishing prototypes for clinical applications; promoting the availability, extensibility, and capacity in non-laboratory applications. (2) To significantly promote the research of ADE in the future. In the discussion we specifically focus on the availability of databases, the transparency of code, the collaboration of research groups, and the imbalanced distribution of training samples.

5.1. The Availability of Databases

Because of the sensitivity of depression data, it is difficult to gain various data for estimating the scale of depression. Hence, the availability of data is a major issue. First, as opposed to the facial expression recognition task, the availability of databases is scarce up until the present day. Given the literature review, one can note that the widely used depression databases are AVEC2013, AVEC2014, DAIC-WOZ. Notably, AVEC2014 is a subset of AVEC2013. Second, there is no multi-modal (i.e., audio, video, text, physiological signals) database to learn comprehensive depression representations for ADE. The existing databases consist of two or three modalities. Though the DAIC database comprises three modalities (audiovisual and text), the organizer has not provided the original videos

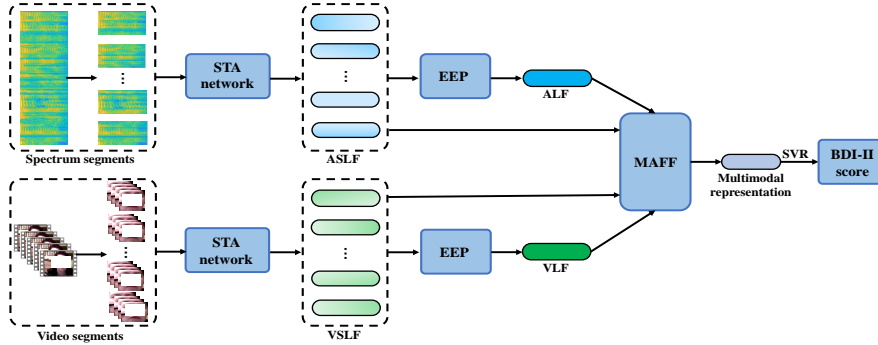


Figure 27: The pipeline of the multi-modal spatiotemporal architecture for deep depression recognition. The approach first inputs spectrogram/video segments into the Spatio-Temporal Attention (STA) network and then uses the features of the last fully connected layers as Audio Segment-Level Feature (ASLF) and Video Segment-Level Feature (VSLF), respectively. Subsequently, Eigen Evolution Pooling (EEP) approach is used for pooling the ASLFs and VSLFs to the ALF and VLF. Lastly, Support Vector Regression (SVR) is used for predicting the BDI-II scores [163].

of DAIC, leading to a certain inconvenience for ADE. Third, the limited size of the data sets limits the research in depression prediction, especially when using DL technologies. For instance, AVEC2013 only contains 50 samples for the training, development, test set, respectively. To address this bottleneck, effective methods to augment the limited amount of annotated data are called for. Fourth, the criteria for data collection should be standardized. At present, different organizers adopt a range of conditions, equipment, and configurations to collect the multi-modal data.

5.2. The Transparency of Data and Algorithms

Despite the significant advances in ADE tasks, there still exist many aspects for improving performance in clinical environments. Nowadays, many researchers from the affective computing community may not share their algorithms using web applications (i.e., github, personal web site). As we know, the DAIC-WOZ dataset has been widely used by many researchers for ADE in the affective computing field. However, the publisher of DAIC-WOZ considers it difficult to provide the data as raw video clips due to the sensitivity of the personal attributes related to mental disorders. Hence, we encourage the all studies to share also the raw data and not only hand-crafted features, or at least to arrange an access to the data in a secure computing environment if sharing is not feasible.

At the very least, researchers should make their code publicly available. Accordingly, different researchers can validate the efficiency of algorithms to build a solid foundation for clinical application. For instance, feature extraction is important for ensemble ADE. However, the current bottleneck is to know which features are suitable for ADE. Given that the principal way of learning features for depression scale prediction is to use deep learning, researchers should design a network that is most suitable for this task. Currently, no commonly accepted standard DL architecture has been defined for ADE.

5.3. The Collaboration of Research Groups

With the significant progress among different disciplines, collaboration with other disciplines is crucial for ADE. For the topic of affective computing, relevant fields include psychology, physiology, computer science, machine learning, etc. Thus, researchers should borrow each other's strengths for promoting the advances of ADE. For audio-based ADE, the deep models only represent the depression scale from audios. Specific to video-based ADE, the deep models capture patterns only from facial expressions. Notably, physiological signals also contain significant information closely related to depression estimation. Accordingly, different researchers should study together to build a multi-modal based DL approaches for clinical application.

5.4. The Imbalanced Distribution of Training Samples

Besides the issues mentioned above, another main issue is the imbalanced distribution of training samples. This issue originates from the fact that the severity of depression is assessed by different discrete numerical values. There are two challenges in modeling the imbalanced data. Firstly, training with imbalanced data may lead to obtaining

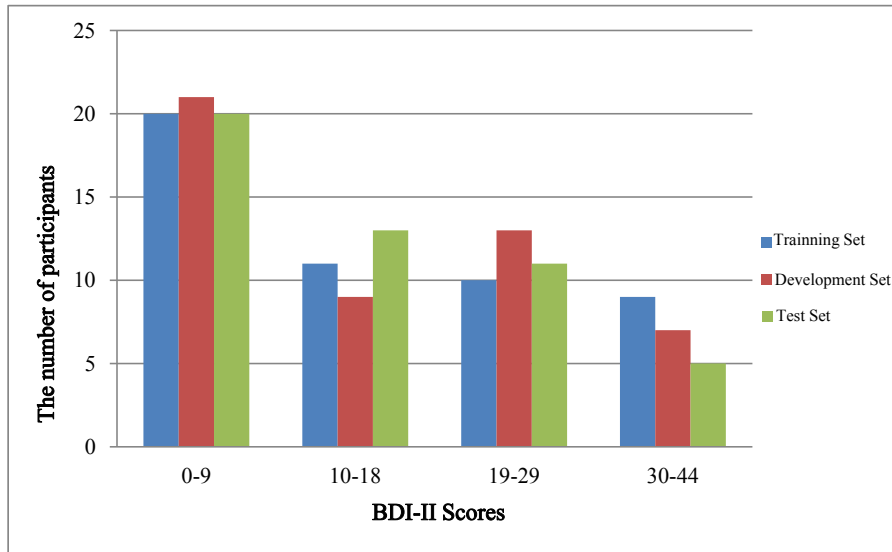


Figure 28: The distribution of BDI-II scores in the training, development, test set of AVEC2014. The range of BDI scores is from 0 to 63 (*no or minimal depression: the range from 0 to 13*), (*Mild: the range from 14 to 19*), (*Moderate: the range from 20 to 28*), (*Severe: the range from 29 to 63*)). The maximum BDI-II score of AVEC2014 is 45. The Y axis is the number of samples in the range of BDI-II.

poor performance of the trained model in the minority category [207]. Secondly, training the models on imbalanced data sets can markedly affect the validation/test set performance. Jeni et al. [208] investigate the effect of skew based on the imbalanced validation set. To address the imbalance, the findings are based on several evaluation metrics, i.e., Accuracy, F1-score, and so on. To demonstrate the imbalance, consider AVEC2014 which is a popular database used by many in the affective computing field. As mentioned in Section 3.2, the BDI-II scores can be divided into four classes according to the depression scales from mild to severity, i.e., 0-9, 10-18, 19-29, 30-44. As illustrated in Fig. 28, the 0-9 class has more and the 30-44 class fewer participants compared to the other classes. Therefore, the database providers should consider the issue of data imbalance to facilitate the training of the shadow or deep models for depression analysis.

6. Conclusions

Along with the progress of deep learning [23], multiple works using DL have also been proposed ADE, with promising performance, establishing the foundations for the clinical application of ADE systems. The present comprehensive survey of the existing ADE methods reviews the topic from several angles while also highlighting numerous issues for further exploration. As a mental disorder, the diagnostics of depression rely on a concerted effort from multiple fields including clinical psychology, affective computing, and computer science. Based on the problems mentioned, developing automatic, objective evaluation systems is valuable for both academic research as well as clinical application. At present, several issues remain to be addressed: 1) the ability to make a distinction between MDD and other depression types [95]; 2) the capacity of learning from few training samples; 3) the ability to extract discriminative features by hand-crafted and deep learning methods; 4) the ability to represent and combine the complementary information from audio and visual cues by fusion approaches.

We conclude by emphasizing the considerable potential in the clinical scenario for assessing the severity of depression. Despite the great progress in recent years, to assist the clinical application, more work should be conducted to collect the additional data, explore a range of methods, and design implement ADE systems for clinical use cases.

In the future, we will resolve the following issues:

1. For the small samples of the training data, on the one hand, we should encourage the data organizer to share the private data samples for ADE. On the other hand, we will try to collect a multi-modal database that includes audio, video, text, physiological signals (i.e., EEG, ECG, etc.). Therefore, different modalities can augment the

data samples for training the ADE models. In addition, we would like to encourage the researchers to share their code on different platforms.

2. To extract the informative features of multi-modal cues, we will consider the attributes of individuals by the DL method. Meanwhile, we will leverage the attributes of data to extract the informative and discriminative features for ADE. In addition, we will collaborate with the researcher from interdisciplinary areas to extract the more informative feature closely related to depression.
3. To learn the complementary patterns between hand-crafted and deep-learned features, robust methods will be designed for ADE. Though deep-learned features have been proven to obtain promising ADE performance, the transnational hand-crafted features should not be ignored for ADE tasks. Therefore, we will deeply study the complementary characteristic between hand-crafted and deep-learned features to model the discriminative architecture for ADE.
4. Multi-modal data are not only augment the size of data to train the modality, but also capture the discriminative patterns for ADE. To improve the performance of the multi-modal ADE, we will consider the complementary patterns among different modality and draw on the experience of researchers from different fields. All in all, we drive the ADE research into clinical application to benefit for those depressed subjects in the future.

7. Acknowledgment

This work is supported by the Scientific Research Program Funded by Shaanxi Provincial Education Department (Program No. 20JG030), the Special Construction Fund for Key Disciplines of Shaanxi Provincial Higher Education, the Shaanxi Higher Education Association Fund for the Prevention and Control of Novel Coronavirus Pneumonia (grant XGH20201), the Shaanxi Provincial Public Scientific Quality Promotion Fund for Emergency Popularization of COVID-19 (grant 2020PSL(Y)040). This work was supported by the Academy of Finland (grants 336033, 315896), Business Finland (grant 884/31/2018), and EU H2020 (grant 101016775). We would also like to thank Dr. Jonathan Cheung-Wai Chan from VUB for proof-reading the manuscript.

References

- [1] C. D. Mathers, D. Loncar, Projections of global mortality and burden of disease from 2002 to 2030, *PLoS medicine* 3 (11) (2006) e442.
- [2] R. C. Kessler, P. Berglund, O. Demler, R. Jin, D. Koretz, K. R. Merikangas, A. J. Rush, E. E. Walters, P. S. Wang, The epidemiology of major depressive disorder: results from the National Comorbidity Survey Replication (NCS-R), *Jama* 289 (23) (2003) 3095–3105.
- [3] K. Hawton, C. C. Comabella, C. Haw, K. Saunders, Risk factors for suicide in individuals with depression: a systematic review, *Journal of affective disorders* 147 (1-3) (2013) 17–28.
- [4] A. McGirr, J. Renaud, M. Seguin, M. Alda, C. Benkelfat, A. Lesage, G. Turecki, An examination of dsm-iv depressive symptoms and risk for suicide completion in major depressive disorder: a psychological autopsy study, *Journal of Affective Disorders* 97 (1-3) (2007) 203–209.
- [5] M. Maj, D. J. Stein, G. Parker, M. Zimmerman, G. Fava, M. D. Hert, K. Demyttenaere, R. McIntyre, T. Widiger, H. Wittchen, The clinical characterization of the adult patient with depression aimed at personalization of management, *World Psychiatry* 19 (2020).
- [6] M. Hamilton, The hamilton rating scale for depression (1986) 143–152.
- [7] J. C. Mundt, P. J. Snyder, M. S. Cannizzaro, K. Chappie, D. S. Geraltz, Voice acoustic measures of depression severity and treatment response collected via interactive voice response (ivr) technology, *Journal of neurolinguistics* 20 (1) (2007) 50–64.
- [8] M. K. Nock, G. Borges, E. J. Bromet, C. B. Cha, R. C. Kessler, S. Lee, Suicide and suicidal behavior, *Epidemiologic reviews* 30 (1) (2008) 133–154.
- [9] T. Sharp, P. J. Cowen, 5-ht and depression: is the glass half-full?, *Current opinion in pharmacology* 11 (1) (2011) 45–51.
- [10] B. Luscher, Q. Shen, N. Sahir, The gabaergic deficit hypothesis of major depressive disorder, *Molecular psychiatry* 16 (4) (2011) 383–406.
- [11] M. O. Poulter, L. Du, I. C. Weaver, M. Palkovits, G. Faludi, Z. Merali, M. Szyf, H. Anisman, Gabaa receptor promoter hypermethylation in suicide brain: implications for the involvement of epigenetic processes, *Biological psychiatry* 64 (8) (2008) 645–652.
- [12] Y. Dwivedi, H. S. Rizavi, R. R. Conley, R. C. Roberts, C. A. Tamminga, G. N. Pandey, Altered gene expression of brain-derived neurotrophic factor and receptor tyrosine kinase b in postmortem brain of suicides, *Archives of general psychiatry* 60 (8) (2003) 804–815.
- [13] J. Gatt, C. Nemeroff, C. Dobson-Stone, R. Paul, R. Bryant, P. Schofield, E. Gordon, A. Kemp, L. Williams, Interactions between bdnf val66met polymorphism and early life stress predict brain and arousal pathways to syndromal depression and anxiety, *Molecular psychiatry* 14 (7) (2009) 681–695.
- [14] J. F. Cohn, T. S. Krueger, I. Matthews, Y. Yang, M. H. Nguyen, M. T. Padilla, F. Zhou, F. De la Torre, Detecting depression from facial actions and vocal prosody, in: 2009 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops, IEEE, 2009, pp. 1–7.

- [15] N. Cummins, J. Epps, M. Breakspear, R. Goecke, An investigation of depressed speech detection: Features and normalization, in: Twelfth Annual Conference of the International Speech Communication Association, 2011.
- [16] J. Joshi, R. Goecke, S. Alghowinem, A. Dhall, M. Wagner, J. Epps, G. Parker, M. Breakspear, Multimodal assistive technologies for depression diagnosis and monitoring, *Journal on Multimodal User Interfaces* 7 (3) (2013) 217–228.
- [17] S. Scherer, G. Stratou, M. Mahmoud, J. Boberg, J. Gratch, A. Rizzo, L.-P. Morency, Automatic behavior descriptors for psychological disorder analysis, in: 2013 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG), IEEE, 2013, pp. 1–8.
- [18] C. Shan, S. Gong, P. W. McOwan, Facial expression recognition based on local binary patterns: A comprehensive study, *Image and vision Computing* 27 (6) (2009) 803–816.
- [19] L. Wen, X. Li, G. Guo, Y. Zhu, Automated depression diagnosis based on facial dynamic analysis and sparse coding, *IEEE Transactions on Information Forensics and Security* 10 (7) (2015) 1432–1441.
- [20] G. Zhao, M. Pietikainen, Dynamic texture recognition using local binary patterns with an application to facial expressions, *IEEE transactions on pattern analysis and machine intelligence* 29 (6) (2007) 915–928.
- [21] Z. Du, W. Li, D. Huang, Y. Wang, Encoding visual behaviors with attentive temporal convolution for depression prediction, in: 2019 14th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2019), IEEE, 2019, pp. 1–7.
- [22] M. Valstar, B. Schuller, K. Smith, F. Eyben, B. Jiang, S. Bilakhia, S. Schnieder, R. Cowie, M. Pantic, AVEC2013: the continuous audio/visual emotion and depression recognition challenge, in: Proceedings of the 3rd ACM international workshop on Audio/visual emotion challenge, 2013, pp. 3–10.
- [23] X. Ma, H. Yang, Q. Chen, D. Huang, Y. Wang, Depaudionet: An efficient deep model for audio based depression classification, in: Proceedings of the 6th international workshop on audio/visual emotion challenge, 2016, pp. 35–42.
- [24] A. Jan, H. Meng, Y. F. B. A. Gaus, F. Zhang, Artificial intelligent system for automatic depression level analysis through visual and vocal expressions, *IEEE Transactions on Cognitive and Developmental Systems* 10 (3) (2017) 668–680.
- [25] S. Song, L. Shen, M. Valstar, Human behaviour-based automatic depression analysis using hand-crafted statistics and deep learned spectral features, in: 2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018), IEEE, 2018, pp. 158–165.
- [26] J. J. B. Uddin Md Azher, L. Young-Koo, Depression level prediction using deep spatiotemporal features and multilayer bilstm, *IEEE Transactions on Affective Computing* (2020).
- [27] M. Al Jazaery, G. Guo, Video-based depression level analysis by encoding deep spatiotemporal features, *IEEE Transactions on Affective Computing* (2018).
- [28] Y. Zhu, Y. Shang, Z. Shao, G. Guo, Automated depression diagnosis based on deep networks to encode facial appearance and dynamics, *IEEE Transactions on Affective Computing* 9 (4) (2017) 578–584.
- [29] W. C. de Melo, E. Granger, A. Hadid, Combining global and local convolutional 3d networks for detecting depression from facial expressions, FG, 2019.
- [30] W. C. de Melo, E. Granger, A. Hadid, Depression detection based on deep distribution learning, ICIP, 2019.
- [31] S. Song, S. Jaiswal, L. Shen, M. Valstar, Spectral representation of behaviour primitives for depression analysis, *IEEE Transactions on Affective Computing* (2020) 1–1.
- [32] M. A. Uddin, J. B. Joolee, Y.-K. Lee, Depression level prediction using deep spatiotemporal features and multilayer bi-lstm, *IEEE Transactions on Affective Computing* (2020).
- [33] L. He, C. Cao, Automated depression analysis using convolutional neural networks from speech, *Journal of Biomedical Informatics* 83 (2018) 103–111.
- [34] N. Cummins, S. Scherer, J. Krajewski, S. Schnieder, J. Epps, T. F. Quatieri, A review of depression and suicide risk assessment using speech analysis, *Speech Communication* 71 (2015) 10–49.
- [35] A. Pampouchidou, P. Simos, K. Marias, F. Meriaudeau, F. Yang, M. Pedititis, M. Tsiknakis, Automatic assessment of depression based on visual cues: A systematic review, *IEEE Transactions on Affective Computing* (2017).
- [36] J. A. Russell, A circumplex model of affect., *Journal of personality and social psychology* 39 (6) (1980) 1161.
- [37] L.-C. Yu, L.-H. Lee, S. Hao, J. Wang, Y. He, J. Hu, K. R. Lai, X. Zhang, Building chinese affective resources in valence-arousal dimensions, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2016, pp. 540–545.
- [38] A. P. Association, et al., Diagnostic and statistical manual of mental disorders (DSM-5®), American Psychiatric Pub, 2013.
- [39] T. Deckersbach, D. D. Dougherty, S. L. Rauch, Functional imaging of mood and anxiety disorders, *Journal of Neuroimaging* 16 (1) (2006) 1–10.
- [40] K. C. Evans, D. D. Dougherty, M. H. Pollack, S. L. Rauch, Using neuroimaging to predict treatment response in mood and anxiety disorders, *Annals of Clinical Psychiatry* 18 (1) (2006) 33–42.
- [41] H. S. Mayberg, A. M. Lozano, V. Voon, H. E. McNeely, D. Seminowicz, C. Hamani, J. M. Schwab, S. H. Kennedy, Deep brain stimulation for treatment-resistant depression, *Neuron* 45 (5) (2005) 651–660.
- [42] A. J. Niemiec, B. J. Lithgow, Alpha-band characteristics in eeg spectrum indicate reliability of frontal brain asymmetry measures in diagnosis of depression., in: 2005 IEEE Engineering in Medicine and Biology 27th Annual Conference, IEEE, 2006, pp. 7517–7520.
- [43] E. J. Nestler, M. Barrot, R. J. DiLeone, A. J. Eisch, S. J. Gold, L. M. Monteggia, Neurobiology of depression, *Neuron* 34 (1) (2002) 13–25.
- [44] R. Cadoret, E. Troughton, L. M. Merchant, A. Whitters, Early life psychosocial events and adult affective symptoms., 1990.
- [45] K. P. Lesch, Gene–environment interaction and the genetics of depression, *Journal of Psychiatry and Neuroscience* 29 (3) (2004) 174.
- [46] R. Cadoret, T. O’Gorman, E. Heywood, E. Troughton, Genetic and environmental factors in major depression., *Journal of affective disorders* 9 2 (1985) 155–64.
- [47] T. A. Brown, P. A. Di Nardo, C. L. Lehman, L. A. Campbell, Reliability of dsm-iv anxiety and mood disorders: implications for the classification of emotional disorders., *Journal of abnormal psychology* 110 (1) (2001) 49.
- [48] J. H. Kamphuis, A. Noordhof, On categorical diagnoses in dsm-v: cutting dimensions at useful points?, *Psychological Assessment* 21 (3) (2009) 294.

- [49] V. Lux, K. Kendler, Deconstructing major depression: a validation study of the dsm-iv symptomatic criteria, *Psychological medicine* 40 (10) (2010) 1679.
- [50] M. A. Oquendo, E. Baca-García, J. J. Mann, J. Giner, Issues for dsm-v: suicidal behavior as a separate diagnosis on a separate axis (2008).
- [51] D. J. Stein, K. A. Phillips, D. Bolton, K. Fulford, J. Z. Sadler, K. S. Kendler, What is a mental/psychiatric disorder? from dsm-iv to dsm-v, *Psychological medicine* 40 (11) (2010) 1759–1765.
- [52] D. Watson, Rethinking the mood and anxiety disorders: a quantitative hierarchical model for dsm-v., *Journal of abnormal psychology* 114 (4) (2005) 522.
- [53] S. D. Østergaard, S. Jensen, P. Bech, The heterogeneity of the depressive syndrome: when numbers get serious., *Acta Psychiatrica Scandinavica* (2011).
- [54] M. JH Balsters, E. J Krahmer, M. GJ Swerts, A. JJM Vingerhoets, Verbal and nonverbal correlates for depression: a review, *Current Psychiatry Reviews* 8 (3) (2012) 227–234.
- [55] W. Chow, M. Doane, J. Sheehan, L. Alphas, H. Le, Le h. economic burden among patients with major depressive disorder: an analysis of healthcare resource use, work productivity, and direct and indirect costs by depression severity, *Am J Manag Care* 16 (2019) e188–e196.
- [56] P. Sobocki, B. Jönsson, J. Angst, C. Rehnberg, Cost of depression in europe., *Journal of Mental Health Policy and Economics* (2006).
- [57] A. J. Mitchell, A. Vaze, S. Rao, Clinical diagnosis of depression in primary care: a meta-analysis, *The Lancet* 374 (9690) (2009) 609–619.
- [58] I. Schumann, A. Schneider, C. Kantert, B. Löwe, K. Linde, Physicians’ attitudes, diagnostic process and barriers regarding depression diagnosis in primary care: a systematic review of qualitative studies, *Family practice* 29 (3) (2012) 255–263.
- [59] R. C. Kessler, E. J. Bromet, The epidemiology of depression across cultures, *Annual review of public health* 34 (2013) 119–138.
- [60] A. T. Beck, R. A. Steer, R. Ball, W. F. Ranieri, Comparison of beck depression inventories-ia and-ii in psychiatric outpatients, *Journal of personality assessment* 67 (3) (1996) 588–597.
- [61] L. Baer, M. A. Blais, *Handbook of clinical rating scales and assessment in psychiatry and mental health*, Springer, 2010.
- [62] D. Maust, M. Cristancho, L. Gray, S. Rushing, C. Tjoa, M. E. Thase, Psychiatric rating scales, in: *Handbook of Clinical Neurology*, Vol. 106, Elsevier, 2012, pp. 227–237.
- [63] R. M. Bagby, A. G. Ryder, D. R. Schuller, M. B. Marshall, The hamilton depression rating scale: has the gold standard become a lead weight?, *American Journal of Psychiatry* 161 (12) (2004) 2163–2177.
- [64] R. D. Gibbons, D. C. Clark, D. J. Kupfer, Exactly what does the hamilton depression rating scale measure?, *Journal of psychiatric research* 27 (3) (1993) 259–273.
- [65] P. Bech, P. Allerup, L. Gram, N. Reisby, R. Rosenberg, O. Jacobsen, A. Nagy, The hamilton depression scale: evaluation of objectivity using logistic models, *Acta Psychiatrica Scandinavica* 63 (3) (1981) 290–299.
- [66] D. Faries, J. Herrera, J. Rayamajhi, D. DeBrotta, M. Demitrack, W. Z. Potter, The responsiveness of the hamilton depression rating scale, *Journal of psychiatric research* 34 (1) (2000) 3–10.
- [67] C. Cusin, H. Yang, A. Yeung, M. Fava, Rating scales for depression, in: *Handbook of clinical rating scales and assessment in psychiatry and mental health*, Springer, 2009, pp. 7–35.
- [68] R. Nuevo, V. Lehtinen, P. M. Reyna-Liberato, J. L. Ayuso-Mateos, Usefulness of the beck depression inventory as a screening method for depression among the general population of finland, *Scandinavian journal of public health* 37 (1) (2009) 28–34.
- [69] L. S. Williams, E. J. Brizendine, L. Plue, T. Bakas, W. Tu, H. Hendrie, K. Kroenke, Performance of the phq-9 as a screening tool for depression after stroke, *stroke* 36 (3) (2005) 635–638.
- [70] P. Pichot, Self-report inventories in the study of depression, in: *New Results in Depression Research*, Springer, 1986, pp. 53–58.
- [71] Y. S. Ben-Porath, Assessing personality and psychopathology with self-report inventories, *Handbook of psychology* (2003) 553–577.
- [72] S. Gilbody, T. Sheldon, A. House, Screening and case-finding instruments for depression: a meta-analysis, *Cmaj* 178 (8) (2008) 997–1003.
- [73] Y. Ren, H. Yang, C. Browning, S. Thomas, M. Liu, Performance of screening tools in detecting major depressive disorder among patients with coronary heart disease: a systematic review, *Medical science monitor: international medical journal of experimental and clinical research* 21 (2015) 646.
- [74] E. Stockings, L. Degenhardt, Y. Y. Lee, C. Mihalopoulos, A. Liu, M. Hobbs, G. Patton, Symptom screening scales for detecting major depressive disorder in children and adolescents: a systematic review and meta-analysis of reliability, validity and diagnostic utility, *Journal of affective disorders* 174 (2015) 447–463.
- [75] A. J. Mitchell, J. C. Coyne, *Screening for depression in clinical practice: an evidence-based guide*, OUP USA, 2009.
- [76] K. Kroenke, R. L. Spitzer, The PHQ-9: A New Depression Diagnostic and Severity Measure, *Psychiatric Annals* 32 (9) (2002) 509–515.
- [77] A. J. Rush, M. H. Trivedi, H. M. Ibrahim, The 16-Item Quick Inventory of Depressive Symptomatology (QIDS), Clinician Rating (QIDS-C), and Self-report (QIDS-SR): A Psychometric Evaluation in Patients with Chronic Major Depression, *Biological Psychiatry* 54 (5) (2003) 573–583.
- [78] S. A. Montgomery, M. Asberg, A new depression scale designed to be sensitive to change., *The British Journal of Psychiatry* 134 (4) (1979) 382–389.
- [79] A. J. Rush, C. M. Gullion, M. R. Basco, R. B. Jarrett, M. H. Trivedi, The inventory of depressive symptomatology (ids): Psychometric properties, *Psychological Medicine* 26 (3) (1996) 477–486.
- [80] W. W. Zung, A self-rating depression scale, *Archives of General Psychiatry* 12 (1) (1965) 63–70.
- [81] H. Ellgring, *Non-verbal communication in depression*, Cambridge University Press, 1989.
- [82] P. H. Waxer, Therapist training in nonverbal communication: I. nonverbal cues for depression., *Journal of clinical psychology* (1974).
- [83] A. Costanza, I. D’Orta, N. Perroud, S. Burkhardt, A. Malafosse, P. Mangin, R. La Harpe, Neurobiology of suicide: do biomarkers exist?, *International journal of legal medicine* 128 (1) (2014) 73–82.
- [84] E. Kraepelin, Manic depressive insanity and paranoia, *The Journal of Nervous and Mental Disease* 53 (4) (1921) 350.
- [85] G. J. Siegle, S. R. Steinhauer, E. S. Friedman, W. S. Thompson, M. E. Thase, Remission prognosis for cognitive therapy for recurrent depression using the pupil: utility and neural correlates, *Biological psychiatry* 69 (8) (2011) 726–733.
- [86] J. S. Silk, R. E. Dahl, N. D. Ryan, E. E. Forbes, D. A. Axelson, B. Birmaher, G. J. Siegle, Pupillary reactivity to emotional information in child and adolescent depression: links to clinical and ecological measures, *American Journal of Psychiatry* 164 (12) (2007) 1873–1880.

- [87] N. P. Jones, G. J. Siegle, D. Mandell, Motivational and emotional influences on cognitive control in depression: A pupillometry study, *Cognitive, Affective, & Behavioral Neuroscience* 15 (2) (2015) 263–275.
- [88] J. Wang, Y. Fan, X. Zhao, N. Chen, Pupillometry in chinese female patients with depression: a pilot study, *International journal of environmental research and public health* 11 (2) (2014) 2236–2243.
- [89] D. Zhou, J. Luo, V. M. Silenzio, Y. Zhou, J. Hu, G. Currier, H. A. Kautz, Tackling mental health by integrating unobtrusive multimodal sensing., in: *AAAI*, 2015, pp. 1401–1409.
- [90] A. Y. Kudinova, K. L. Burkhouse, G. Siegle, M. Owens, M. L. Woody, B. E. Gibb, Pupillary reactivity to negative stimuli prospectively predicts recurrence of major depressive disorder in women, *Psychophysiology* 53 (12) (2016) 1836–1842.
- [91] M. Li, S. Lu, G. Wang, L. Feng, B. Fu, N. Zhong, Alleviated negative rather than positive attentional bias in patients with depression in remission: an eye-tracking study, *Journal of International Medical Research* 44 (5) (2016) 1072–1086.
- [92] R. B. Price, D. Rosen, G. J. Siegle, C. D. Ladouceur, K. Tang, K. B. Allen, N. D. Ryan, R. E. Dahl, E. E. Forbes, J. S. Silk, From anxious youth to depressed adolescents: Prospective prediction of 2-year depression symptoms via attentional bias measures., *Journal of Abnormal Psychology* 125 (2) (2016) 267.
- [93] G. Stratou, S. Scherer, J. Gratch, L.-P. Morency, Automatic nonverbal behavior indicators of depression and ptsd: Exploring gender differences, in: *2013 Humaine Association Conference on Affective Computing and Intelligent Interaction, IEEE*, 2013, pp. 147–152.
- [94] S. Ghosh, M. Chatterjee, L.-P. Morency, A multimodal context-based approach for distress assessment, in: *Proceedings of the 16th International Conference on Multimodal Interaction*, 2014, pp. 240–246.
- [95] G. Stratou, S. Scherer, J. Gratch, L.-P. Morency, Automatic nonverbal behavior indicators of depression and ptsd: the effect of gender, *Journal on Multimodal User Interfaces* 9 (1) (2015) 17–29.
- [96] Z. Yu, S. Scherer, D. Devault, J. Gratch, G. Stratou, L.-P. Morency, J. Cassell, Multimodal prediction of psychological disorders: Learning verbal and nonverbal commonalities in adjacency pairs, in: *SemDial 2013 DialDam: Proceedings of the 17th Workshop on the Semantics and Pragmatics of Dialogue*, 2013, pp. 160–169.
- [97] L.-P. Morency, G. Stratou, D. DeVault, A. Hartholt, M. Lhomme, G. M. Lucas, F. Morbini, K. Georgila, S. Scherer, J. Gratch, et al., Sinsensei demonstration: A perceptive virtual human interviewer for healthcare applications., in: *AAAI*, 2015, pp. 4307–4308.
- [98] G. M. Lucas, J. Gratch, S. Scherer, J. Boberg, G. Stratou, Towards an affective interface for assessment of psychological distress, in: *2015 International Conference on Affective Computing and Intelligent Interaction (ACII), IEEE*, 2015, pp. 539–545.
- [99] S. Scherer, G. Stratou, L.-P. Morency, Audiovisual behavior descriptors for depression assessment, in: *Proceedings of the 15th ACM International conference on multimodal interaction*, 2013, pp. 135–140.
- [100] J. Gratch, R. Artstein, G. M. Lucas, G. Stratou, S. Scherer, A. Nazarian, R. Wood, J. Boberg, D. DeVault, S. Marsella, et al., The distress analysis interview corpus of human and computer interviews., in: *LREC*, 2014, pp. 3123–3128.
- [101] S. Scherer, G. Stratou, G. Lucas, M. Mahmoud, J. Boberg, J. Gratch, L.-P. Morency, et al., Automatic audiovisual behavior descriptors for psychological disorder analysis, *Image and Vision Computing* 32 (10) (2014) 648–658.
- [102] S. Alghowinem, R. Goecke, M. Wagner, G. Parker, M. Breakspear, Eye movement analysis for depression detection, in: *2013 IEEE International Conference on Image Processing*, 2013, pp. 4220–4224.
- [103] S. Alghowinem, R. Goecke, J. F. Cohn, M. Wagner, G. Parker, M. Breakspear, Cross-cultural detection of depression from nonverbal behaviour, in: *2015 11th IEEE International conference and workshops on automatic face and gesture recognition (FG)*, Vol. 1, *IEEE*, 2015, pp. 1–8.
- [104] R. Gupta, N. Malandrakis, B. Xiao, T. Guha, M. Van Segbroeck, M. Black, A. Potamianos, S. Narayanan, Multimodal prediction of affective dimensions and depression in human-computer interactions, in: *Proceedings of the 4th International Workshop on Audio/Visual Emotion Challenge*, 2014, pp. 33–40.
- [105] P. Ekman, W. Friesen, J. Hager, Facial action coding system: The manual. research nexus, *Network Information Research Corp.*, Salt Lake City, UT 1 (2002) 8.
- [106] G. McIntyre, R. Göcke, M. Hyett, M. Green, M. Breakspear, An approach for automatically measuring facial activity in depressed subjects, in: *2009 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops, IEEE*, 2009, pp. 1–8.
- [107] G. J. McIntyre, et al., The computer analysis of facial expressions: on the example of depression and anxiety (2010).
- [108] J. F. Cohn, Social signal processing in depression, in: *Proceedings of the 2nd international workshop on Social signal processing*, 2010, pp. 1–2.
- [109] G. McIntyre, R. Goecke, M. Breakspear, G. Parker, Facial response to video content in depression, in: *MMCogEmS Workshop: Inferring Cognitive and Emotional States from Multimodal Measures*, 13th International Conference on Multimodal Interaction ICMI2011, Alicante, Spain, 2011.
- [110] J. M. Girard, J. F. Cohn, M. H. Mahoor, S. M. Mavadati, Z. Hammal, D. P. Rosenwald, Nonverbal social withdrawal in depression: Evidence from manual and automatic analyses, *Image and vision computing* 32 (10) (2014) 641–647.
- [111] J. M. Girard, J. F. Cohn, M. H. Mahoor, S. Mavadati, D. P. Rosenwald, Social risk and depression: Evidence from manual and automatic facial expression analysis, in: *2013 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, *IEEE*, 2013, pp. 1–8.
- [112] J. F. Cohn, Beyond group differences: specificity of nonverbal behavior and interpersonal communication to depression severity., in: *AVEC@ ACM Multimedia*, Citeseer, 2013, pp. 1–2.
- [113] J. R. Williamson, T. F. Quatieri, B. S. Helfer, G. Ciccarelli, D. D. Mehta, Vocal and facial biomarkers of depression based on motor incoordination and timing, in: *Proceedings of the 4th International Workshop on Audio/Visual Emotion Challenge*, 2014, pp. 65–72.
- [114] M. K. Mandal, A. Awasthi, Understanding facial expressions in communication: Cross-cultural and multidisciplinary perspectives, *Springer*, 2014.
- [115] T.-H. Yang, C.-H. Wu, K.-Y. Huang, M.-H. Su, Coupled hmm-based multimodal fusion for mood disorder detection through elicited audio-visual signals, *Journal of Ambient Intelligence and Humanized Computing* 8 (6) (2017) 895–906.
- [116] M. Valstar, J. Gratch, B. Schuller, F. Ringeval, D. Lalanne, M. Torres Torres, S. Scherer, G. Stratou, R. Cowie, M. Pantic, Avec 2016: Depression, mood, and emotion recognition workshop and challenge, in: *Proceedings of the 6th international workshop on audio/visual*

- emotion challenge, 2016, pp. 3–10.
- [117] L. Yang, D. Jiang, L. He, E. Pei, M. C. Oveneke, H. Sahli, Decision tree based depression classification from audio video and language information, in: Proceedings of the 6th international workshop on audio/visual emotion challenge, 2016, pp. 89–96.
- [118] L. Yang, H. Sahli, X. Xia, E. Pei, M. C. Oveneke, D. Jiang, Hybrid depression classification and estimation from audio video and text information, in: Proceedings of the 7th Annual Workshop on Audio/Visual Emotion Challenge, 2017, pp. 45–51.
- [119] L. Yang, D. Jiang, H. Sahli, Integrating deep and shallow models for multi-modal depression analysis—hybrid architectures, *IEEE Transactions on Affective Computing* (2018).
- [120] S. Alghowinem, R. Goecke, M. Wagner, G. Parker, M. Breakspear, Head pose and movement analysis as an indicator of depression, in: 2013 Humaine Association Conference on Affective Computing and Intelligent Interaction, 2013, pp. 283–288.
- [121] J. Joshi, Depression analysis: a multimodal approach, in: Proceedings of the 14th ACM international conference on Multimodal interaction, 2012, pp. 321–324.
- [122] J. Joshi, An automated framework for depression analysis, in: 2013 Humaine Association Conference on Affective Computing and Intelligent Interaction, IEEE, 2013, pp. 630–635.
- [123] J. Joshi, R. Goecke, G. Parker, M. Breakspear, Can body expressions contribute to automatic depression analysis?, in: 2013 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG), IEEE, 2013, pp. 1–7.
- [124] J. Joshi, A. Dhall, R. Goecke, M. Breakspear, G. Parker, Neural-net classification for spatio-temporal descriptor based depression analysis, in: Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012), IEEE, 2012, pp. 2634–2638.
- [125] J. Joshi, A. Dhall, R. Goecke, J. F. Cohn, Relative body parts movement for automatic depression analysis, in: 2013 Humaine Association Conference on Affective Computing and Intelligent Interaction, IEEE, 2013, pp. 492–497.
- [126] B. Hosseinifard, M. H. Moradi, R. Rostami, Classifying depression patients and normal subjects using machine learning techniques and nonlinear features from eeg signal, *Computer methods and programs in biomedicine* 109 (3) (2013) 339–345.
- [127] R. Adomi, A. Gatti, A. Brugnera, K. Sakatani, A. Compare, Could fnirs promote neuroscience approach in clinical psychology?, *Frontiers in psychology* 7 (2016) 456.
- [128] C. S. Ho, L. J. Lim, A. Lim, N. H. Chan, R. Tan, S. Lee, R. Ho, Diagnostic and predictive applications of functional near-infrared spectroscopy for major depressive disorder: A systematic review, *Frontiers in Psychiatry* 11 (2020) 378.
- [129] T. Suto, M. Fukuda, M. Ito, T. Uehara, M. Mikuni, Multichannel near-infrared spectroscopy in depression and schizophrenia: cognitive brain activation study, *Biological psychiatry* 55 (5) (2004) 501–511.
- [130] S. Scherer, G. M. Lucas, J. Gratch, A. S. Rizzo, L.-P. Morency, Self-reported symptoms of depression and ptsd are associated with reduced vowel space in screening interviews, *IEEE Transactions on Affective Computing* 7 (1) (2015) 59–73.
- [131] J. T. Becker, F. Boiler, O. L. Lopez, J. Saxton, K. L. McGonigle, The natural history of cognitive decline in alzheimer’s disease, *Archives of Neurology* 51 (6) (1994) 585–594.
- [132] H. Stassen, S. Kuny, D. Hell, The speech analysis approach to determining onset of improvement under antidepressants, *European Neuropsychopharmacology* 8 (4) (1998) 303–310.
- [133] D. J. France, R. G. Shiavi, S. Silverman, M. Wilkes, Acoustical properties of speech as indicators of depression and suicidal risk, *IEEE Transactions on Biomedical Engineering* 47 (7) (2000) 829–837.
- [134] M. Alpert, E. R. Pouget, R. R. Silva, Reflections of depression in acoustic measures of the patient’s speech, *Journal of Affective Disorders* 66 (1) (2001) 59–69.
- [135] E. Moore, M. Clements, J. Peifer, L. Weisser, Comparing objective feature statistics of speech for classifying clinical depression, in: *Engineering in Medicine and Biology Society, 2004. IEMBS’04. 26th Annual International Conference of the IEEE*, Vol. 1, IEEE, IEEE, San Francisco, CA, USA, 2004, pp. 17–20.
- [136] T. Yingthawornskuk, H. K. Keskinpala, D. France, D. M. Wilkes, R. G. Shiavi, R. M. Salomon, Objective estimation of suicidal risk using vocal output characteristics, in: *Ninth International Conference on Spoken Language Processing, ISCA, Pittsburgh, Pennsylvania, 2006*, pp. 649–652.
- [137] N. C. Maddage, R. Senaratne, L.-S. A. Low, M. Lech, N. Allen, Video-based detection of the clinical depression in adolescents, in: *Engineering in Medicine and Biology Society, 2009. EMBC 2009. Annual International Conference of the IEEE*, IEEE, IEEE, Minneapolis, MN, USA, 2009, pp. 3723–3726.
- [138] S. Alghowinem, R. Goecke, M. Wagner, J. Epps, M. Breakspear, G. Parker, From joyous to clinically depressed: Mood detection using spontaneous speech, in: *Proceedings of the Twenty-Fifth International Florida Artificial Intelligence Research Society Conference, AAAI, Marco Island, Florida, 2012*, pp. 141–146.
- [139] K. Ooi, L. Low, M. Lech, N. Allen, Prediction of clinical depression in adolescents using facial image analysis, in: *WIAMIS 2011: 12th International Workshop on Image Analysis for Multimedia Interactive Services, WIAMIS, WIAMIS, Delft, The Netherlands, 2011*, pp. 1–4.
- [140] J. C. Mundt, A. P. Vogel, D. E. Feltner, W. R. Lenderking, Vocal acoustic biomarkers of depression severity and treatment response, *Biological Psychiatry* 72 (7) (2012) 580–587.
- [141] M. Valstar, B. Schuller, K. Smith, T. Almaev, F. Eyben, J. Krajewski, R. Cowie, M. Pantic, AVEC 2014: 3D dimensional affect and depression recognition challenge, in: *Proceedings of the 4th International Workshop on Audio/Visual Emotion Challenge, ACM, ACM, Orlando, Florida, USA, 2014*, pp. 3–10.
- [142] G. G. J. Chen, Visualizations for mental health topic models (2014) 1–54.
- [143] K.-Y. Huang, C.-H. Wu, Y.-T. Kuo, F.-L. Jang, Unipolar depression vs. bipolar disorder: An elicitation-based approach to short-term detection of mood disorder., in: *INTERSPEECH, 2016*, pp. 1452–1456.
- [144] H. Dibeklioglu, Z. Hammal, J. F. Cohn, Dynamic multimodal measurement of depression severity using deep autoencoding, *IEEE Journal of Biomedical and Health Informatics* 22 (2) (2018) 525–536.
- [145] E. Çiftçi, H. Kaya, H. Güleç, A. A. Salah, The turkish audio-visual bipolar disorder corpus, in: *2018 First Asian Conference on Affective Computing and Intelligent Interaction (ACII Asia)*, IEEE, 2018, pp. 1–6.
- [146] H. Cai, Y. Gao, S. Sun, N. Li, F. Tian, H. Xiao, J. Li, Z. Yang, X. Li, Q. Zhao, et al., Modma dataset: a multi-model open dataset for mental-disorder analysis, *arXiv preprint arXiv:2002.09283* (2020).

- [147] S. Alghowinem, R. Goecke, M. Wagner, J. Epps, M. Hyett, G. Parker, M. Breakspear, Multimodal depression detection: fusion analysis of paralinguistic, head pose and eye gaze behaviors, *IEEE Transactions on Affective Computing* 9 (4) (2016) 478–490.
- [148] E. Moore, M. Clements, J. Peifer, L. Weisser, Analysis of prosodic variation in speech for clinical depression, in: *Proceedings of the 25th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (IEEE Cat. No. 03CH37439)*, Vol. 3, IEEE, 2003, pp. 2925–2928.
- [149] H. Meng, D. Huang, H. Wang, H. Yang, M. Ai-Shuraifi, Y. Wang, Depression recognition based on dynamic facial and vocal expression features using partial least square regression, in: *Proceedings of the 3rd ACM international workshop on Audio/visual emotion challenge*, 2013, pp. 21–30.
- [150] N. Cummins, J. Joshi, A. Dhall, V. Sethu, R. Goecke, J. Epps, Diagnosis of depression by behavioural signals: a multimodal approach, in: *Proceedings of the 3rd ACM international workshop on Audio/visual emotion challenge*, 2013, pp. 11–20.
- [151] K. Ooi, Early prediction of clinical depression in adolescents using single-channel and multi-channel classification approach (2014).
- [152] M. Sidorov, W. Minker, Emotion recognition and depression diagnosis by acoustic and visual features: A multimodal approach, in: *Proceedings of the 4th International Workshop on Audio/Visual Emotion Challenge*, 2014, pp. 81–86.
- [153] M. Kächele, M. Schels, F. Schwenker, Inferring depression and affect from application dependent meta knowledge, in: *Proceedings of the 4th International Workshop on Audio/Visual Emotion Challenge*, 2014, pp. 41–48.
- [154] I. T. Meftah, N. Le Thanh, C. B. Amar, Detecting depression using multimodal approach of emotion recognition, in: *2012 IEEE International Conference on Complex Systems (ICCS)*, IEEE, 2012, pp. 1–6.
- [155] U. R. Acharya, S. L. Oh, Y. Hagiwara, J. H. Tan, H. Adeli, D. P. Subha, Automated eeg-based screening of depression using deep convolutional neural network, *Computer Methods and Programs in Biomedicine* 161 (2018) 103–113.
- [156] A. Zandvakili, N. S. Philip, S. R. Jones, A. R. Tyrka, B. D. Greenberg, L. L. Carpenter, Use of machine learning in predicting clinical response to transcranial magnetic stimulation in comorbid posttraumatic stress disorder and major depression: A resting state electroencephalography study, *Journal of Affective Disorders* 252 (2019) 47–54.
- [157] D. Kan, P. Lee, Decrease alpha waves in depression: An electroencephalogram (eeg) study, in: *2015 International Conference on BioSignal Analysis, Processing and Systems (ICBAPS)*, IEEE, 2015, pp. 156–161.
- [158] X. Zhang, J. Shen, Z. ud Din, J. Liu, G. Wang, B. Hu, Multimodal depression detection: Fusion of electroencephalography and paralinguistic behaviors using a novel strategy for classifier ensemble, *IEEE Journal of Biomedical and Health Informatics* 23 (6) (2019) 2265–2275.
- [159] H. Cai, X. Zhang, Y. Zhang, Z. Wang, B. Hu, A case-based reasoning model for depression based on three-electrode eeg data, *IEEE Transactions on Affective Computing* (2018) 1–1.
- [160] D. Zhi, X. Ma, L. Lv, Q. Ke, Y. Yang, X. Yang, M. Pan, S. Qi, R. Jiang, Y. Du, et al., Abnormal dynamic functional network connectivity and graph theoretical analysis in major depressive disorder, in: *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, IEEE, 2018, pp. 558–561.
- [161] J. J. Maller, S. S. George, R. P. Viswanathan, P. B. Fitzgerald, P. Junor, Using thermographic cameras to investigate eye temperature and clinical severity in depression, *Journal of biomedical optics* 21 (2) (2016) 026001.
- [162] M. Hamilton, Arating scalefordepression. *journalofneurology, Neurosurgery, and Psychiatry* 23 (1960) 5642.
- [163] M. Niu, J. Tao, B. Liu, J. Huang, Z. Lian, Multimodal spatiotemporal representation for automatic depression level detection, *IEEE Transactions on Affective Computing* (2020).
- [164] P. Viola, M. J. Jones, Robust real-time face detection, *International journal of computer vision* 57 (2) (2004) 137–154.
- [165] T. Baltrušaitis, P. Robinson, L.-P. Morency, Openface: an open source facial behavior analysis toolkit, in: *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 2016, pp. 1–10.
- [166] G. Littlewort, J. Whitehill, T. Wu, I. Fasel, M. Frank, J. Movellan, M. Bartlett, The computer expression recognition toolbox (cert), in: *Face and gesture 2011*, IEEE, 2011, pp. 298–305.
- [167] M. Cannizzaro, B. Harel, N. Reilly, P. Chappell, P. J. Snyder, Voice acoustical measurement of the severity of major depression, *Brain and cognition* 56 (1) (2004) 30–35.
- [168] E. Moore II, M. A. Clements, J. W. Peifer, L. Weisser, Critical analysis of the impact of glottal features in the classification of clinical depression in speech, *IEEE transactions on biomedical engineering* 55 (1) (2007) 96–107.
- [169] B. J. Shannon, K. K. Paliwal, A comparative study of filter bank spacing for speech recognition, in: *Microelectronic engineering research conference*, Vol. 41, Citeseer, 2003, pp. 310–12.
- [170] F. Eyben, B. Schuller, opensmile:) the munich open-source large-scale multimedia feature extractor, *ACM SIGMultimedia Records* 6 (4) (2015) 4–13.
- [171] L. Yang, D. Jiang, H. Sahli, Feature augmenting networks for improving depression severity estimation from speech signals, *IEEE Access* 8 (2020) 24033–24045.
- [172] M. Niu, B. Liu, J. Tao, Q. Li, A time-frequency channel attention and vectorization network for automatic depression level prediction, *Neurocomputing* (2021). doi:<https://doi.org/10.1016/j.neucom.2021.04.056>. URL <https://www.sciencedirect.com/science/article/pii/S0925231221005981>
- [173] Y. Dong, X. Yang, A hierarchical depression detection model based on vocal and emotional cues, *Neurocomputing* 441 (2021) 279–290.
- [174] D. Yi, Z. Lei, S. Liao, S. Z. Li, Learning face representation from scratch, *arXiv preprint arXiv:1411.7923* (2014).
- [175] X. Zhou, K. Jin, Y. Shang, G. Guo, Visually interpretable representation learning for depression recognition from facial images, *IEEE Transactions on Affective Computing* 11 (3) (2020) 542–552. doi:10.1109/TAFFC.2018.2828819.
- [176] W. C. de Meto, E. Granger, M. B. Lopez, Encoding temporal information for automatic depression recognition from facial analysis, in: *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2020, pp. 1080–1084.
- [177] L. He, J. C.-W. Chan, Z. Wang, Automatic depression recognition using cnn with attention mechanism from videos, *Neurocomputing* 422 165–175.
- [178] M. Niu, J. Tao, B. Liu, C. Fan, Automatic depression level detection via lp-norm pooling., in: *INTERSPEECH*, 2019, pp. 4559–4563.
- [179] Z. Zhao, Q. Li, N. Cummins, B. Liu, H. Wang, J. Tao, B. W. Schuller, Hybrid network feature extraction for depression assessment from speech, *Proc. Interspeech 2020* (2020) 4956–4960.

- [180] T. Al Hanai, M. M. Ghassemi, J. R. Glass, Detecting depression with audio/text sequence modeling of interviews., in: Interspeech, 2018, pp. 1716–1720.
- [181] W. C. de Melo, E. Granger, A. Hadid, A deep multiscale spatiotemporal network for assessing depression from facial dynamics, *IEEE Transactions on Affective Computing* (2020).
- [182] L. He, C. Guo, P. Tiwari, H. M. Pandey, W. Dang, Intelligent system for depression scale estimation with facial expressions and case study in industrial intelligence, *International Journal of Intelligent Systems* (2021).
- [183] W. Carneiro de Melo, E. Granger, M. Bordallo Lopez, Mdn: A deep maximization-differentiation network for spatio-temporal depression detection, *IEEE Transactions on Affective Computing* (2021) 1–1doi:10.1109/TAFFC.2021.3072579.
- [184] L. Yang, D. Jiang, W. Han, H. Sahli, Dnn and dnn based multi-modal depression recognition, in: 2017 Seventh International Conference on Affective Computing and Intelligent Interaction (ACII), IEEE, 2017, pp. 484–489.
- [185] Z. Zhao, Z. Bao, Z. Zhang, J. Deng, N. Cummins, H. Wang, J. Tao, B. Schuller, Automatic assessment of depression from speech via a hierarchical attention transfer network and attention autoencoders, *IEEE Journal of Selected Topics in Signal Processing* 14 (2) (2020) 423–434. doi:10.1109/JSTSP.2019.2955012.
- [186] L. Yang, D. Jiang, X. Xia, E. Pei, M. C. Oveneke, H. Sahli, Multimodal measurement of depression using deep learning models, in: Proceedings of the 7th Annual Workshop on Audio/Visual Emotion Challenge, 2017, pp. 53–59.
- [187] S. Yin, C. Liang, H. Ding, S. Wang, A multi-modal hierarchical recurrent neural network for depression detection, in: Proceedings of the 9th International on Audio/Visual Emotion Challenge and Workshop, 2019, pp. 65–71.
- [188] M. Rodrigues Makiuchi, T. Warnita, K. Uto, K. Shinoda, Multimodal fusion of bert-cnn and gated cnn representations for depression detection, in: Proceedings of the 9th International on Audio/Visual Emotion Challenge and Workshop, 2019, pp. 55–63.
- [189] W. Fan, Z. He, X. Xing, B. Cai, W. Lu, Multi-modality depression detection via multi-scale temporal dilated cnns, in: Proceedings of the 9th International on Audio/Visual Emotion Challenge and Workshop, 2019, pp. 73–80.
- [190] L. Zhang, J. Driscoll, X. Chen, R. Hosseini Ghomi, Evaluating acoustic and linguistic features of detecting depression sub-challenge dataset, in: Proceedings of the 9th International on Audio/Visual Emotion Challenge and Workshop, 2019, pp. 47–53.
- [191] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* 30 (2017) 5998–6008.
- [192] K. Hara, H. Kataoka, Y. Satoh, Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet?, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, 2018, pp. 6546–6555.
- [193] L. Yang, Y. Li, H. Chen, D. Jiang, M. C. Oveneke, H. Sahli, Bipolar disorder recognition with histogram features of arousal and body gestures, in: Proceedings of the 2018 on Audio/Visual Emotion Challenge and Workshop, 2018, pp. 15–21.
- [194] Z. Du, W. Li, D. Huang, Y. Wang, Bipolar disorder recognition via multi-scale discriminative audio temporal representation, in: Proceedings of the 2018 on Audio/Visual Emotion Challenge and Workshop, 2018, pp. 23–30.
- [195] Z. S. Syed, K. Sidorov, D. Marshall, Automated screening for bipolar disorder from audio/visual modalities, in: Proceedings of the 2018 on Audio/Visual Emotion Challenge and Workshop, 2018, pp. 39–45.
- [196] X. Xing, B. Cai, Y. Zhao, S. Li, Z. He, W. Fan, Multi-modality hierarchical recall based on gbdt for bipolar disorder classification, in: Proceedings of the 2018 on Audio/Visual Emotion Challenge and Workshop, 2018, pp. 31–37.
- [197] Z. Huang, J. Epps, D. Joachim, Investigation of speech landmark patterns for depression detection, *IEEE Transactions on Affective Computing* (2019) 1–1doi:10.1109/TAFFC.2019.2944380.
- [198] K. Anis, H. Zakia, D. Mohamed, C. Jeffrey, Detecting depression severity by interpretable representations of motion dynamics, in: 2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018), IEEE, 2018, pp. 739–745.
- [199] S. M. Alghowinem, T. Gedeon, R. Goecke, J. Cohn, G. Parker, Interpretation of depression detection models via feature selection methods, *IEEE Transaction on Affective Computing* (01) (2020) 1–1.
- [200] S. Jayawardena, J. Epps, E. Ambikairajah, Ordinal logistic regression with partial proportional odds for depression prediction, *IEEE Transactions on Affective Computing* (2020).
- [201] L. He, D. Jiang, H. Sahli, Automatic depression analysis using dynamic facial appearance descriptor and dirichlet process fisher encoding, *IEEE Transactions on Multimedia* 21 (6) (2018) 1476–1486.
- [202] S. Jayawardena, J. Epps, E. Ambikairajah, Evaluation measures for depression prediction and affective computing, in: ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2019, pp. 6610–6614. doi:10.1109/ICASSP.2019.8682956.
- [203] S. Amiriparian, A. Awad, M. Gerczuk, L. Stappen, A. Baird, S. Ottl, B. Schuller, Audio-based recognition of bipolar disorder utilising capsule networks, in: 2019 International Joint Conference on Neural Networks (IJCNN), IEEE, 2019, pp. 1–7.
- [204] Y. Li, L. Yang, H. Chen, D. Jiang, H. Sahli, Audio visual multimodal classification of bipolar disorder episodes, in: 2019 8th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos (ACIIW), IEEE, 2019, pp. 115–120.
- [205] Z. Ren, J. Han, N. Cummins, Q. Kong, M. D. Plumbley, B. W. Schuller, Multi-instance learning for bipolar disorder diagnosis using weakly labelled speech data, in: Proceedings of the 9th International Conference on Digital Public Health, 2019, pp. 79–83.
- [206] N. Abaei, H. Al Osman, A hybrid model for bipolar disorder classification from visual information, in: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2020, pp. 4107–4111.
- [207] H. He, E. A. Garcia, Learning from imbalanced data, *IEEE Transactions on knowledge and data engineering* 21 (9) (2009) 1263–1284.
- [208] L. A. Jeni, J. F. Cohn, F. De La Torre, Facing imbalanced data—recommendations for the use of performance metrics, in: 2013 Humaine association conference on affective computing and intelligent interaction, IEEE, 2013, pp. 245–251.

Table 3: Summary of the Audiovisual databases which have been adopted in the reviewed works for the last 20 years. Abbreviations: DPRD – Depressed, SCDL – Suicidal, NTRL – Neutral, not depressed or suicidal, M – Number of males, F – Number of Females DSM - Diagnostic and Statistical Manual of Mental Disorders, HAMID - Hamilton Rating Scale for Depression, BDI -Beck Depression Inventory, QIDS - Quick Inventory of Depressive Symptomology, PHQ-9 - Patient Health Questionnaire. Note: where DSM is present as a clinical score for all depressed patients in corpus to meet criteria for Major Depressive Disorder

Database	Modality	Subjects	Annotation	Ground Truth	Public /Private
①: DementiaBank [131](1994)	A+V+T	226	HAMD	Clinical Assessment	Public
②: – [132](1998)	A	43	HAMD (DPRS = HAMD ≥ 10)	–	Private
③: – [133](2000)	A	115	DSM-IV, BDI (DPRS = BDI > 20)	Clinical Assessment	Private
④: – [134](2001)	A	41	DSM-III-R, HAMD (DPRS = HAMD ≥ 18)	Clinical Assessment	Private
⑤: – [135](2004)	A	33	DSM-IV	Clinical Assessment	Private
⑥: – [136](2006)	A	32	BDI (DPRD = BDI > 20)	Interviews with a Therapist	Private
⑦: – [14](2009)	A	57	DSM-IV, HAMD (DPRD = HAMD ≥ 15)	–	Private
⑧: ORI [137](2009)	V	139	Manual annotation	–	Private
⑨: BlackDog [138](2009)	A+V	80	DSM-IV, HAMD > 15	Clinical Assessment	Private
⑩: ORYGEN [139](2011)	V	191	Manual annotation	Clinical Assessment	Private
⑪: – [140] (2012)	Audio	165	HAMD, DSM-IV, QIDS-C, QIDS (QIDS-SR)	Clinical assessment	Private
⑫: AVEC2013 [22] (2013)	A+V	292	BDI-II	Self-report	Public
⑬: AVEC2014 [141] (2014)	A+V	292	BDI-II	Self-report	Public
⑭: Crisis Text Line [142] (2014)	T	–	Manual annotation	–	Private
⑮: DAIC-WoZ [100] (2014)	A+V+T	110	PHQ-9 (DPRD = PHQ-9 > 10)	Self-report	Public
⑯: Rochester [89] (2015)	V	27	Manual annotation	Self-report	Private
⑰: CHI-MEI [143](2016)	V	53	DSSS, HAM-D	Clinical Assessment	Private
⑱: Pittsburgh [144] (2018)	A+V	57	DSM-IV, HAMD > 15	Clinical Assessment	Public
⑲: BD [145] (2018)	A+V	46	DSM-V	Clinical Assessment	Public
⑳: MODMA [146] (2020)	A+EEG	EEG-128(53), EEG-3(55), A(55)	PHQ-9	Clinical Assessment	Public

Table 4: Performance summary of reviewed approaches for depression recognition based on static images of the most widely assessed databases. Note that the list results of DAIC-WOZ and E-DAIC contain text features for ADE.

Modality	Datasets	Methods	Network Type	Preprocessing	Test	
					RMSE	MAE
Audio	AVEC2013	He et al. 2018 [33]	DCNN	DFT	10.00	8.20
		Niu et al. 2019 [178]	DCNN, LSTM	-	9.79	7.48
		Zhao et al. 2020 [179]	DCNN	-	9.65	7.38
		Niu et al. 2021 [172]	DCNN	-	8.32	6.26
		Dong et al. 2021 [173]	DCNN	-	8.73	7.31
	AVEC2014	He et al. 2018 [33]	DCNN	DFT	9.99	8.19
		Niu et al. 2019 [178]	DCNN, LSTM	-	9.66	8.02
		Zhao et al. 2020 [179]	DCNN	-	9.57	7.94
		Niu et al. 2021 [172]	DCNN	-	9.25	7.49
		Dong et al. 2021 [173]	DCNN	-	8.82	6.79
	DAIC-WOZ	Ma et al. 2016 [23]	DCNN, LSTM	-	-	-
		Alhanai et al. 2018 [180]	LSTM	-	6.50	5.13
		Yang et al. 2020 [171]	DCGAN	-	5.52	4.63
Video (Static Images)	AVEC2013	Zhu et al. 2017 [28]	DCNN	Dlib	9.82	7.58
		Zhou et al. 2018 [175]	DCNN	Dlib	8.28	6.20
		Melo et al. 2019 [30]	RetNet-50	MTCNN	8.25	6.30
		Melo et al. 2020 [176]	DCNN	MTCNN	7.97	5.96
		Sone et al. 2020 [31]	DCNN	OpenFace	8.10	6.16
		He et al. 2020 [177]	DCNN	OpenFace	8.39	6.59
	AVEC2014	Zhou et al. 2018 [175]	DCNN	Dlib	8.39	6.21
		Zhu et al. 2017 [28]	DCNN	Dlib	9.55	7.47
		He et al. 2020 [177]	DCNN	OpenFace	8.30	6.51
		Melo et al. 2019 [30]	RetNet-50	MTCNN	8.23	6.15
		Melo et al. 2020 [176]	DCNN	MTCNN	7.94	6.20
		Sone et al. 2020 [31]	DCNN	OpenFace	7.15	5.95
		Video (Image Sequences)	AVEC2013	Jazaery et al. 2018 [27]	RNN, C3D	OpenFace
Melo et al. 2019 [29]	C3D			MTCNN	8.26	6.40
Azher et al. 2020 [26]	DCNN, LSTM			-	8.93	7.04
Melo et al. 2020 [181]	3D-CNN			MTCNN	7.90	5.98
He et al. 2021 [182]	DCNN			OpenFace	8.46	6.83
Melo et al. 2021 [183]	MDN			MTCNN	7.55	6.24
AVEC2014	Jazaery et al. 2018 [27]		RNN, C3D	OpenFace	9.20	7.22
	Melo et al. 2019 [29]		C3D	MTCNN	8.31	6.59
	Azher et al. 2020 [26]		DCNN, LSTM	-	8.78	6.86
	Melo et al. 2020 [181]		3D-CNN	MTCNN	7.61	5.82
	He et al. 2021 [182]		DCNN	OpenFace	8.42	6.78
	Melo et al. 2021 [183]		MDN	MTCNN	7.65	6.06
DAIC-WOZ	Song et al. 2018 [25]		DCNN	-	5.84	4.37
	Du et al. 2019 [21]	DCNN	-	5.78	4.61	
Multi-modal	AVEC2013	Niu et al. 2020 [163]	DCNN	Dlib, STFT	8.16	6.14
	AVEC2014	Niu et al. 2020 [163]	DCNN	Dlib, STFT	7.03	5.21
		Jan et al. 2018 [24]	DCNN	-	7.43	6.14
	DAIC-WOZ	Yang et al. 2017 [184]	DCNN, DNN	-	6.34	5.38
		Yang et al. 2018 [119]	DCNN, DNN	-	6.34	5.39
		Zhao et al. 2020 [185]	BLSTM	-	5.51	4.20
		Yang et al. 2017 [186]	DCNN	-	5.97	5.16
		Yang et al. 2017 [118]	DCNN	-	5.40	4.35
		E-DAIC	Shi et al. 2019 [187]	LSTM	-	5.50
	Makiuchi et al. 2019 [188]		DCNN, LSTM	-	6.11	-
	Fan et al. 2019 [189]		DCNN	-	5.91	4.39
Zhang et al. 2019 [190]	Random Forest		-	6.85	5.84	

Table 5: Performance summary of reviewed approaches for multi-modal depression recognition on the BD databases.

Datasets	Methods	Network Type	Preprocessing	Dev.	Test
				UAR	
BD	Yang et al. 2018 [193]	DCNN	Openpose		57.41%
	Xing et al. 2018 [196]	GBDT	–		57.41%
	Du et al. 2018 [194]	LSTM	–	65.1%	–
	Shahin et al. 2019 [203]	DCNN	–	–	45.5%
	Sun et al. 2021 [203]	DCNN	–	93.12%	–
	Li et al. 2019 [204]	DCNN	–	74.5%	–
	Ren et al. 2019 [205]	DCNN, DNN	–	–	57.41%
	Abaei et al. 2020 [206]	DCNN, LSTM	–	60.67%	–