

EMOTION DETECTION USING SPEECH RECOGNIZATION AND FACIAL EXPRESSIONS

K.Neha Nandini¹, G.Varshitha², K.Sai Preethi³, N.Rikhitha⁴., P.Deepthi⁵

1 Assistant Professor, Department of CSE., Malla Reddy College of Engineering for Women., Maisammaguda., Medchal., TS, India (✉ nehanandini.kella@gmail.com)

2, 3, 4, 5 B.Tech CSE, (19RG1A05574, 19RG1A0591, 19RG1A05A3, 19RG1A05A6), Malla Reddy College of Engineering for Women., Maisammaguda., Medchal., TS, India

Abstract

A lot of progress has been made in the field of automated facial emotion recognition (FER) in recent years. FER has been implemented in human-centered computing and the emerging field of emotional AI (EAI), both of which aim to improve human-machine interactions. The goal of EAI research is to improve computers' ability to read and interpret human emotions and actions in a wide variety of contexts. In recent years, neural networks have undergone tremendous evolution, leading to the development of new designs to tackle ever more challenging tasks. Deep learning has been the most influential technique in this regard. In this post, we'll look at the state of the art in AI-based automatic emotion identification utilizing cutting-edge deep learning models. We demonstrate that deep learning-based FER can work in tandem with other models that employ architecture-related approaches, such as databases, to provide very precise outcomes.

Introduction

Facial expressions may be easily recognized thanks to the major and unique characteristics of the human face. Definition of FER - When a person's outward expression betrays their true feelings. It finds widespread usage in computer vision, digital image processing, and artificial intelligence, as well as HCI applications including face image processing, facial video surveillance, and facial animation. The challenging problem of automatic facial expression identification has attracted the attention of a growing number of scientists in recent years. The feature extraction process is vital in FER. A.L. and coworkers[1] that 55% of all communication is conveyed via facial expressions, whereas just 38% is through verbal and spoken communication.

There are two main approaches to FER system design. Some methods use a series of pictures, beginning with a neutral expression and progressing to the most extreme ones, as a training phase. However, systems that just employ a single face picture to make an emotion recognition determination tend to underperform when compared to state-of-the-art methods [2,3]. A FER system may be classified in more than one way, depending on the attributes it uses for recognition, in addition to the approach type it models. categories. The first group of characteristics is gleaned from the orientation of the face organs and the feel of the skin. Geometric features are the second sort of feature; they store data about the face's different locations and points and may be used to evaluate

either a single picture or a series of photographs by tracking the face's motion across the frame. Extracting geometric characteristics from a face might begin with identifying key facial landmarks. Landmarks are prominent facial features that may be used to learn more about a person's identity from observation alone. Numerous more experiments have been conducted on facial landmark identification, however these are outside the purview of this paper. In order to locate these landmarks, this study makes use of the Python library dlib [4]. There are two distinct aspects of AI that are involved in the automated recognition of human emotions and psychology. Researchers in the fields of psychology and artificial intelligence are focusing on answering the challenge of how to recognize human emotions. The tone and auditory modifications [5] that are widely available, for example [6] and from which a rapid mood evaluation may be obtained, and other sources [7] are examples of the vocal and nonverbal sensors that create subjects like mood and accent. Mehrian's research [8] found that the senses accounted for 55% of information (emotional and linguistic) and the other 7% had some kind of physical component that was not defined. Many scientists are keen to study facial expressions since they are the first reliable indicator of an individual's emotional state.

To begin, we'll be expanding our extraction features. improvements to other aspects, adding features to an existing representation might be beneficial. Each coded movement in the Facial Action Coding System (FACS) is assumed to require at least one facial muscle, as highlighted by Ekman and Friesen [9]. When it comes to coding head motions, Ekman and Friesen were the first to notice how FACS facial movement is implemented in FACS facial AUs. (among several volunteers of various sexes and/or genetic backgrounds.

Purpose of the Investigation 1.1. There are several universal and fundamental facial expressions that may be used to convey a wide range of human emotions. If an algorithm can be built to recognize, extract, and evaluate these facial expressions in real time, then it will be feasible to automatically identify emotions in still photographs and movies. Expressions on one's face may convey a wealth of information about one's state of mind and goals in a social context. They're crucial to how we communicate with one another as social beings. Facial expression processing greatly benefits from the *additional context provided by seeing faces in their natural environments. Understanding and expressing empathy in interpersonal interactions is*

essential. Automatic emotion recognition has always been a hotly debated topic in psychology. Because of this, a lot of development has occurred in this field. Words, gestures, and expressions on the face and body are only few of the ways in which we communicate our emotions.

expressions. Therefore, it is essential for human-machine communication that both parties can read and understand expressions of emotion.

1.1. 1.1. This Study's Original Contribution

In this study, we survey the state of the art in emotion recognition using several architectures for identifying a wide variety of expressive styles from face signals alone. The most up-to-date findings from 2016-2021 are published, and analyses of the most common issues and current efforts to fixing them are provided as well. How it works is as follows. Section 2 introduces the prototypical facial expressions and other foundational kinds used to define face expressions, such as FACS. In Section 3, we see the framework of our system for recognizing people's facial expressions and moods. Practical applications of current FER findings are discussed in Section 4. The difficulties encountered by FER in the region are briefly discussed in Section 5, followed by some predictions for the future. Basic forms of emotion recognition are discussed in Section 7, after which links to public databases used in FER tasks are provided in Section 6. In Section 8, we get a quick rundown of how deep learning can be used to recognize emotions in people's faces. The ninth section is a discussion and comparison of FER. In the conclusion, we look forward to what the future may contain. **2. The Most Common Expression Categories**

There are two primary approaches to consider when characterizing facial expressions. A Coding System for Facial Expressions 2.1 The FACS [11] is able to detect subtle changes in face features. See Figure 1; the muscles of facial expression are the 1. frontalis, 2. orbicularis oculi, 3. zygomaticus major, 4. risorius, 5. platysma, and 6. depressor anguli oris. This widely used method in psychology is based on the observation of a human observer and consists of 44 action units connected to the tightening of groups of facial muscles in order to detect facial emotions. Figure 2 also depicts a few of the action components. Skilled personnel typically classify and label FACS manually, looking at slow-motion video footage of facial muscle contractions. There have been several recent attempts to automate this process [13]. The system's ability to record complex facial expressions, among other things, makes it vulnerable to the potential challenge posed by its reliance on descriptive data labels rather than inferential data labels. The FACS data must be turned into a system capable of estimating emotions before it can be used. The Emotional Facial Action System (EFACS) [14] is a paradigm for this sort of behavior.

2.2. Prototypic Emotional Expressions. Most FER systems go right to defining prototype expressions rather than defining facial details. The human universal facial expression of emotion set [15] is the most widely used collection of prototype facial emotion expressions, and it covers six types of basic emotions. It's the go-to library of expression

templates for usage in communication. These base concepts are used due to their universal applicability (Figure 3). This suggests that these feelings are shared by all humans and may be seen in many contexts [17]. One may express emotions such as fear, anger, excitement, sorrow, distaste, surprise, or indifference. This system may either be used as a standard classifier to determine what emotion the person in the picture is feeling, or as a probabilistic estimator to determine how likely it is that the individual is indeed feeling that emotion. It acts as a fuzzy classifier in the second scenario.

Third, the framework for emotion recognition and face detection

FER may function both independently and as a plug-in module for other types of face recognition software. That's why it's smart to look at the system's bigger picture. As can be seen in Figure 3, the system is comprised of four main parts. The input material is analyzed by the face detection module to determine whether a face is present.

If the source material is a video, just the most important frames will undergo facial recognition, while the rest will be tracked for any changes. This is done to increase the robustness of the system as a whole. However, face alignment is similar to face detection in that it pinpoints the exact location of the detected face. During this stage, we are tasked with recognizing many aspects of a person's face. Then, the picture isophotometric properties, like as brightness and contrast, were modified using a method called geometric normalization. Labels such as gender, identity, and expression are then classified using feature extraction. Depending on the circumstances, the extracted feature may be sent into a classifier or compared to training data.

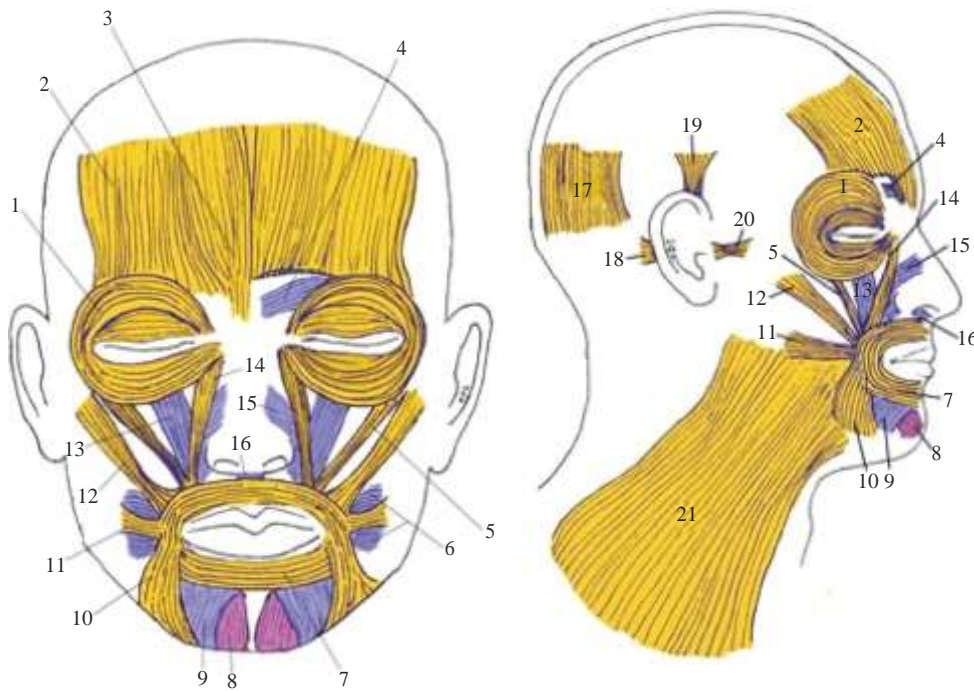
In Sections 3.1, 3.2, and 3.3, we go into the nuts and bolts of face detection, face alignment, and feature extraction, respectively. Face Recognition System. Detecting faces is the first step in the face recognition process [14], and it's crucial to the success of the whole system. Movie characters may often be identified by their movements, expressions, and skin tones. The face is the primary focus of many of the most successful techniques for changing one's look [18]. The difficulties in modeling 3D structures like faces may be avoided by these approaches. However, the face/nonface boundary is sometimes somewhat murky, hence 3D variations are required for emotion recognition in faces. Since the 1990s, various solutions [19] have been presented to address this problem. Kenli and Ai [20] developed a method for detecting anomalies using Eigen decomposition. They use a variety of "eigenfaces" in addition to a generic face. The study's authors [21] made a distinction between this and the work of Sung and Poggio, who looked only at 'eigenfaces. However, the probability of nonfaces was calculated using Bayes' rule. Using neural networks, Rowley et al. [22] distinguished between face and non-face pictures, while Osuna et al. [23] trained a Kernel support vector machine to do the same. The bootstrap method was used to retrain the SVM, with promising outcomes. Also, Schneiderman and Kanade [24] utilized AdaBoost to create a classifier that takes into account the wavelet structure of an image. Therefore, a considerable amount of processing time is needed to run the technique. By exchanging the wavelets for Haar features [26], Viola and Jones [25] were able to solve this

issue. When compared to wavelets, Haar features required less processing time. A real-time, front-view face recognition system has now been shown for the first time [27].

Viola's framework has been the subject of several suggested improvements. Lienhart and coworkers [28] rotated the Haar characteristics in the plane. To handle out-of-plane rotation, Li et al. [29,30] proposed utilizing a detector pyramid, which may also be used for multiview face recognition. The facial recognition methods Eigenface

and AdaBoost were presented. The Eigenface method is the simplest, while the AdaBoost method is the most effective, for detecting human faces. AdaBoost might potentially be used to identify people by their faces.

Congruency of the Faces 2.12. When used with face localization, face alignment, which includes the recognition of facial feature points, may provide better results. Figure 3 presents a side-by-side evaluation of several face recognition algorithms and facial alignment techniques.



Upper Face Action Units					
Lid Droop	Slit	Eyes Closed	Squint	Blink	Wink
Lower Face Action Units					

Figure 1: Muscles of facial expression [10].

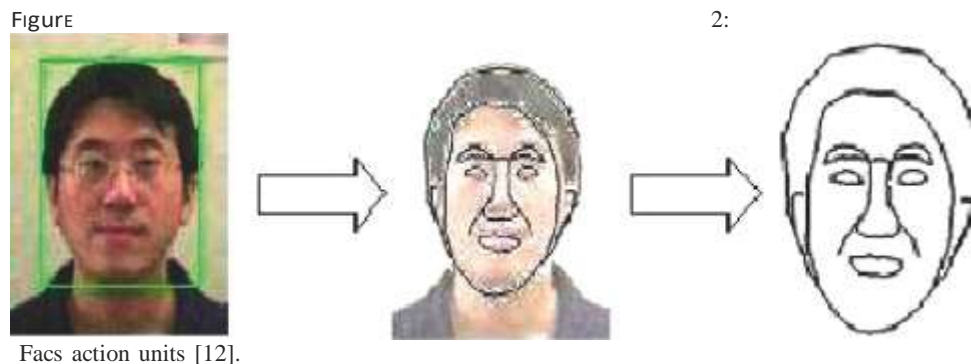


Figure 3: Face detection and alignment processes [16].

Facial alignment is precise down to the pixel level, whereas face detection evaluates larger areas of an image. There have been several proposed solutions to this issue since the 1990s. Histograms were used by Gu et al. [31] to determine where on a photograph the lips and eyes were located. Images processed with Gabor filters were used by Marian and colleagues [32] to locate the medial cleft and pupils. The Active Shape Model [33] is the most efficient curve fitting algorithm known, while many other approaches have been explored.

The Active Shape Model (ASM) was introduced by Cootes et al. [33-35] specifically for facial photos. Since then, significant improvements have been made to the ASM's durability, velocity, and precision. By combining Gabor filters with ASM, Li et al. [36, 37] developed the Direct Appearance Model, which has now been confirmed by further research. Different Writers

Improved ASM for local searches with the use of 2D local textures has been published in [38].

Feature Extraction, Version 2.12. Pixel data is converted into more abstract representations of the face in the picture, such as the face's texture, color, motion, contours, and spatial arrangement, using the feature extraction method. In subsequent classification processes, this information will be utilized to aid in the discovery of trends. It is common practice throughout the feature extraction process to reduce the dimensionality of the input space. High-quality knowledge retention is essential. This approach requires steadiness and discrimination while keeping a cool head. Multiple distinguishing features are used for facial recognition [39].

Eigenface coefficients have been used as features as recently as [40], while an eigenface extension known as Ten- sorface has also showed promising results. The picture of the face is broken down into "shape" and "texture" in the Active Appearance Model [41]. While the shape vector represents the facial outlines, the texture

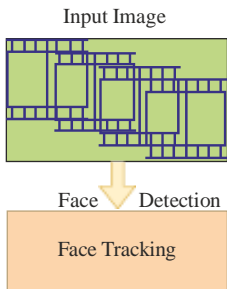
vector describes the "shape-free" textures. Using a two-dimensional mesh, Matsuno et al. [42] extracted features using Potential Net. All of the aforementioned techniques are considered holistic since they consider the whole picture while analyzing an image. Local features are another kind of feature that is hyper-focused on localized areas. In other cases, as those of Colmenarez et al. [43], local features may be utilized directly as picture subwindows.

The 2.16 iteration made use of nine ancillary panes arranged around the characteristics of the face. Another common wavelet filter that has been utilized with some success in terms of vision in the primary visual cortex is the Gabor Filter [44, 45]. Yin and Wei [46] have also employed rudimentary topographical details to symbolize human faces. Instead of explicitly defining the traits, Yu and Bhanu [47] used an evolutionary method to produce them automatically. Video-based FER also includes the dynamic variation of expression. The suggested Geometric Deformation Feature in [48] may translate landmark nodes in a geometric sense. Aleksic and Katsaggelos' [1] Facial Animation Parameters are derived from the Active Shape Model.

The Organization of Feelings (2.16). There have been several attempts to solve the challenge of automated expression recognition using different classifiers. In order to classify face expressions, Matsuno et al. [49] examined the cutoff value of normalized Euclidean distance between features. Bayesian recognition [43] is another approach that uses maximum likelihood to identify facial emotions. The literature also mentions Higher-Order Singular Value Decomposition [52], Locally Linear Embedding [50], and Fisher discrimination analysis [51], among others [1,53,54]. Currently, the most effective solutions to the automated expression recognition issue are neural networks [58–61] and support vector machines [55–57]

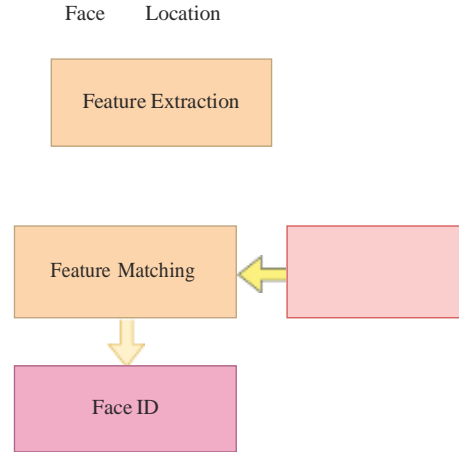
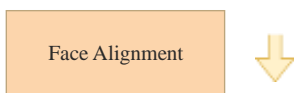
5. Using FER in the Present Moment Only lately

have limitations been identified in the emotional systems [62]. Studies have shown that emotional mechanisms, which may be beneficial or harmful depending on their presence or absence, take priority over cognitive processes in the brain [63]. Bad moods give birth to bad ideas, which limit creativity while trying to find answers to problems, and ultimately bring you deeper into difficulties. Blood flow to the brain exhibits diverse patterns for each emotional state, including anger, sorrow, fear, and happiness [64]. Many studies have shown that happy feelings like pleasure, acceptance, trust, and satisfaction help students learn, whereas anxious feelings like fear, anger, and shame may hinder it. result in impairments in learning and the process as a whole. The effects of anxiety and depression on memory are multifaceted. Some of the ways these conditions manifest themselves are via stress, which is exacerbated by hopelessness and feeds on itself to produce emotions of rage and terror. Intelligent feedback may help students overcome their lack of desire and learn effectively while they are struggling. The latter requires a computer that can read its users' emotions and adapt accordingly, allowing for the development of weak areas of knowledge.



in good standing, manage their interests, and offer them with relevant information and timely response [8]. The whole processing flow for face recognition is shown in Figure 4.

Embodied conversational agents (or other virtual agents that are capable of communicating both verbally and nonverbally, such as animated graphical characters) that can convey emotions or other sentiments and provide information to them using body language are one way to enable students to provide affective and intelligent feedback in an e-learning system, for example [65]. Having a computer that communicates with the user but ignores their input is far more efficient. Implementing such systems is very difficult given our limited capacity to detect human emotions and actions at the moment. Andrew G. Howard et al. [66] created a group of effective convolutional neural models for use in MobileNets. A



kind of productive He also contributed to the development of the MobileNets convolution model class. The innovative depth separable convolution used in the creation of the neural architecture consists of components, rather than linked layers, of depthwise convolutions. The first layer is a depthwise convolution that consists of two independent convolutions, and the two layers may be separated. Compared to conventional convolution, it's nearly 10 times more productive. It does nothing except apply lowpass and highpass filters to the signal, however, so it's not that useful. They accomplished this by doing operations on top of each other, such as pointwise convolution (e.g., which adds up the results of pointwise convolutions) and implementing 1x pointwise convolution. They take on two more hyperparameters in an effort to boost efficiency. As a result, the network may be made more efficient by decreasing the cost of each layer and increasing the network's breadth by a multiplicative factor and resolution by a nonuniform. The model builder may use a variety of hyperparameters to find the optimal balance between accuracy and complexity for their specific needs. This model shows off its many features and capabilities with a wide range of illustrations, from measuring the dimensions of an item to analyzing a person's facial expressions. Emotions are elusive, making them challenging to identify.

naive and so prone to making mistakes; yet, they may often be uncovered using a wide range of tools. Happiness, anger, sorrow, fear, disgust, surprise, and contempt are only some of the eleven fundamental emotions that Ekman thinks may be deduced from seeing people's faces [67]. This received a boost just after the new century.

Figure 4: Face recognition processing flow.

The positive results from tests with face recognition and other forms of audio-visual media have provided a significant boost for the development of research into automated affect recognition. The phrase "to judge by the look on their face" refers to the common belief that recognizing emotions is as simple as using facial expression to seek for patterns that indicate whether or not a person is empathic. Using the Facial Action Coding System (FACS), we may categorize the wide variety of facial actions, including facial gestures, into various AUs, each of which has its own distinct entities, and, ultimately, emotions.

According to Bartlett and Mattivi [68], Bartic and colback's analysis of the literature on emotion recognition is thorough and accurate. Geometrical characteristics (such how the eyes are shaped or the earlobes are positioned) of the eyebrows) assess the dimensions of the face, including the nose and lips. To categorize the identified face into an emotional state, empirical methods conduct feature analysis using a variety of machine learning algorithms. Using the data from the programs created by [69], it is now feasible to detect emotions using facial expressions. The Face Reader is a face analysis software application that accurately detects six fundamental emotions. Researchers in the field of face recognition have made significant strides by exploiting local features in a persons database, as detailed in the study cited in [70]. Both a person's degree of stress and their level of both the emotional investment and the physical stamina required, which is underappreciated.

An individual's attentiveness and emotional condition may be inferred from their gestures and postures to a substantial degree. There is a dearth of study in this area, however these themes have not been the subject of considerable analysis. The examination of the user's past and present interactions with the online data using the aforementioned sources might provide evidence of the user's current cognitive state [71]. Emotional state and entertainment, as well as the research of many types of influencing elements, are crucial to the construction of an affective guiding system, which in turn yields appropriate feedback. To better reflect their students' person- alities and adapt to their emotional states, robotic tutors are anticipated to provide revamped virtual lectures. Researchers in [72] argue that all four feelings—frustration, boredom, motivation, and confidence—are equally important in a computer instructor, and they analyze the many forms of feedback they've collected to determine their relative importance. There are [73] pieces written by the writers that deal with fundamental feelings including fear, grief, and happiness. Before presenting oneself with emotion and voice in a second kind of ECAs called "reactive empathy," one engages in a "expansive empathy" ECA.

The authors provide a difficulty scale developed by Philipp et al. [9] that illustrates several methods of expansion, despite the fact that previous research has shown that automated expansion is the most challenging. When employing this method, it is crucial to take into account the subject's head pose, skin condition, and/or age, which may vary depending on the subject's position, the time of day, and the amount of available light, as well as the problem of occlusion caused by the scarf or other source of illumination. Facial features can be extracted using a number of different methods, including geometric features like LBP [74] and Gaborlet unit activity, and texture features like the Generalized Local Binary Pattern Classification (LBC) and the Directionalized Gabor (GDA) [75]. Since its widespread use in recent years, especially with the help of convolutional neural networks and recurrent neural networks, emotion recognition has become a very effective method. To aid in the creation of deep architectures, many neural networks have been designed; all of them provide respectable results [76]

Current Problems/Challenges in Face Detection and Emotion Recognition

In this article, state-of-the-art techniques for decoding human facial expressions are examined. Face detection and alignment, normalization of the facial picture, extraction of important features, and classification are all crucial steps in developing a facial emotion recognition system. Most systems today still carry out these steps in a sequential and separate manner. AsThis section will thus first examine the challenges of emotion recognition, before moving on to an analysis of how these processes have been dealt with in different research.

Recognizing individual facial characteristics and deducing an individual's emotional state, however, are challenging tasks. The human face is not uniform, and there are additional constraints associated with lighting, shadows, facial location, and orientation issues in different settings, all of which contribute to the difficulty [77]. While humans have an innate talent for reading and understanding facial expressions and emotions, computers continue to struggle with basic tasks like learning to distinguish between different faces. Multilayer perceptron (MLP) neural networks and support vector machines are only two examples of the deep learning approaches that have been investigated as a family to improve the accuracy and speed of fundamental machine learning classification methods. Human behavior analysis works best when used in a variety of settings. It's possible that deep learning algorithms may provide the necessary resilience and scalability when applied to novel sorts of data.

In the sections that follow, we'll examine the most crucial

problems that arise when trying to automate the identification of facial expressions. Obtaining task-representative data, coping with occlusions, modeling dynamics, and overcoming ground truth gathering issues are all significant challenges in this context. Figure 5 depicts the processes used by common FER methods, including the detection of a face region and facial landmarks in input images, the extraction of spatial and temporal features from the face components and landmarks, and the determination of a facial expression using pretrained pattern classifiers based on one of the facial categories (face images are taken from the CK + dataset [78]).

Emotion Recognition in Human Faces: Eight Databases

As face recognition technology improves and becomes more widespread annually, facial database sizes have grown dramatically [79]. Model enhancement or training necessitates a database of examples of the type needed for recognition, as well as class labels for them, and this database must grow in size as the number of examples used in the model increases [80]. Emotional detection, for instance, might be applied in a number of contexts, from basic human-robot cooperation [81] to the detection of depressive symptoms [82].

An alternate variant is one in which the top and bottom half are aligned but cropped differently; this is because the algorithm often takes image/portrait datasets that are evenly illuminated and fixed in po- sition, as seen in the top section. The NIST mugshot database [83] provides a

clear, grayscale alternative for finding picture IDs of 1573 people on a neutral backdrop, which may be compared to

the pixelated versions. On the other hand, the writers have to go a

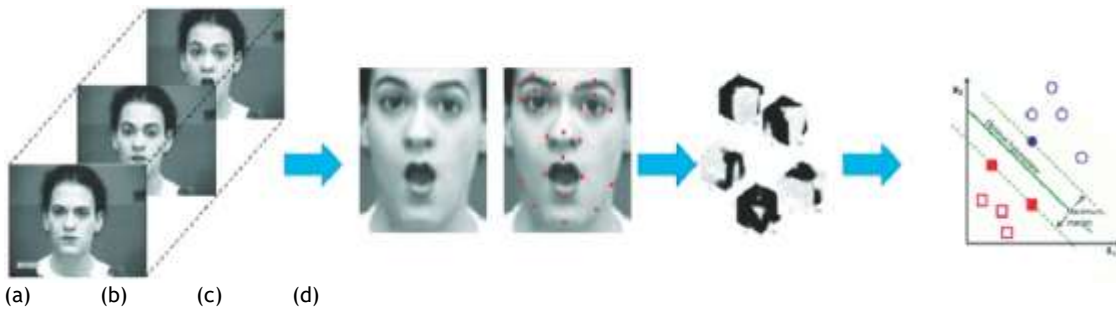


Figure 5: Conventional FER method [28]. (a) Input images. (b) Face detection and landmark detection. (c) Feature extraction. (d) FE classification.

out into the actual world to get a feel for how light conditions and occlusions interact in the context of real-world circumstances, so that you may better understand the situations [84]. The M2VTS database includes the faces of 37 participants in a range of rotated and illuminated locations; this technique [85] was used to effortlessly rotate the subjects and study the effects of different lighting on their look. The value of a database is determined by the kinds of feelings it stores. The six main emotion categories defined by Ekman are used by several databases, including CK, MMI, eINTERFACE, and NVIE. Efforts have been made to categorize or include generic good and negative emotions in several datasets, including the SMO, AAI, and ISL meeting corpus. The CSC corpus database is one example of an attempt to rate dishonesty and honesty. BU-3DFE, BU-4DFE, Bosphorus, and BP4D are some of the most well-known 3D datasets. Six-expression posed datasets are available in both BU-3DFE and BU-4DFE, with the latter having greater resolution. Bosphorus aims to improve the situation by adding more expressions to the avatar's face, whereas BP4D is the only one among the four that employs induced emotions rather than posed ones. The primary advantage of deep learning is that it exposes neural networks to other databases, hence enabling them to expand by incorporating a plethora of fresh inputs, examples, face expressions, and ongoing expression modifications..

8. Multimodal Emotion Recognition: Facial Expression, Spoken Expression, and More

Facial expression recognition, verbal emotion recognition, and multimodal emotion recognition including visual representations are the three primary kinds of emotion recognition methods discussed here. We have also discussed the potential applications of these techniques. *Recognizing Emotion from a Person's Face 7.1. When it comes to nonverbal communication, facial expressions are crucial. Several fields benefit greatly from facial expression recognition technology, including healthcare and human-computer interaction. According to Mehrabian, just 7% of communication takes place in writing, 38% in speech, and 55% in body language. The six primary emotions identified by Ekman and his colleagues [86] are joy, sorrow, surprise, fear, and rage.*

He proved that individuals of all backgrounds have these feelings. Reference: Feldman et al. [32] suggest that valence and arousal, two orthogonal variables, may be used to express emotions. He realized that people express their emotions in different ways. of addition, there is a wide range of replies when individuals are asked to report their emotions on a regular basis [87]. Both positive and negative valences and arousal levels are possible [88]. Information would be classified in this study according to changes in valence and arousal levels. The manual extraction of facial expressions was originally created by researchers via the creation of algorithms for extracted functions such the Gabor wavelet, the Weber Local Descriptor (WLD), the Local Binary Pattern (LBP), and multifeature fusion. These features are vulnerable to uneven coverage of topics, which may cause significant texture data loss in the source picture. The study of facial expressions via the use of deep neural network models is currently the most talked-about topic in the field of facial recognition. In addition, FER has several practical uses in everyday life, such as intelligent security, deception detection, and smart medicine. In [89], the authors explored the use of DBN, deep CNN, and long short-term memory (LSTM) [90], as well as their combination, to create face expression recognition models.

Expression Analysis in Spoken Language 7.1. Human-computer interface systems rely heavily on speech recognition. They will express themselves verbally and nonverbally. Speech recognition algorithms are often used to determine emotional state [91]. Early work on emotion detection in speech focused on extracting artificial characteristics from human speech in order to classify it. Using a set of continuous speech characteristics based on pitch, amplitude, and spectral tilt, Liscombe et al. (2003) investigated the association between different emotions and these sounds. Numerous algorithms have been developed throughout time [92] to identify the range of human emotions conveyed in speech. Various machine learning methods, such as support vector machines, hidden Markov models, and Gaussian mixed models, have been presented. Voice recognition is only one of several speech domains where deep learning has been used successfully [93]. Emotion detection using convolutional neural networks has also been attempted;

these studies demonstrate the superiority of bi-directional multimodal emotion recognitional RNNs (Bi-LSTM) in extracting crucial speech components, thereby enhancing speech recognition. performance [1]. Figure 6 illustrates the end-to-end "SpeechEmotion Recognition" system.

Multimodal Feeling Identification System 7. Many studies still make use of multi-modal emotion processing. Emotion research would benefit from this expansion since it would allow for the incorporation of additional study modalities (video, audio, sensor data, etc.). The research employs a number of strategies and methods in order to complete its objective. Big data, semantics, and deep learning are all used by many of them. Emotions are difficult to identify because they are complex psychophysiological processes that occur nonverbally. There is substantial evidence that multimodal learning is superior than unimodal learning [94]. Recognized faces provide a rich source of visual data that may form the backbone of a neural network for multimodal emotion identification. Their strategy was motivated by the successful entries to the EmotiW competition in 2013 and 2014. This method was developed by Chen et al. [95] as a solution to the problem of multimodal emotion detection (MEC 2016). In order to ascertain the mood of the video's protagonist, this method retrieves multimodal characteristics. When it comes to identifying feelings, the facial CNN feature is the most accurate. In a previous study [96], we used classic and deep convolutional neural network (DCNN) techniques to recover a number of characteristics. This method yields very encouraging results on test data. We detail the methods that were utilized to create the MEC 2017 team submissions for the 2017 Multimodal Emotion Recognition Challenge at Beijing Normal University. A Dempster-Shafer theory fusion approach was offered for combining several prediction results based on the recovered features, which included an autoencoder (AE), a CNN, a dense SIFT, and an audio feature. Figure 7 depicts the architecture for a NN capable of recognizing emotions across several modalities.

In addition, studies have attempted to integrate information from many modalities, including vocal tone, facial emotions, and eye gaze.

combinations of text, physiological signals, and other channels [97]. Emotion recognition accuracy is presently being improved by using this method. Emotion detection findings may be generated via a multimodal fusion model by combining several types of physiological data. Thanks to recent developments in DL architectures, deep learning may now be used for multimodal emotion identification. Deep learning encompasses a wide variety of approaches, including LSTM [55], SVM [98], deep convolutional neural networks [76], and deep belief networks [77].

Facial Expression Analysis Using Deep Neural Networks

Since deep learning algorithms allow automatic feature learning quickly, they have lately emerged as a potential alternative to conventional feature design methods. Research in deep learning may pave the way for more accurate representations and novel algorithms to learn such representations from unlabeled data. The development of high-powered GPU processors has made it possible to do high-performance numerical GPU computing has made these methods computationally

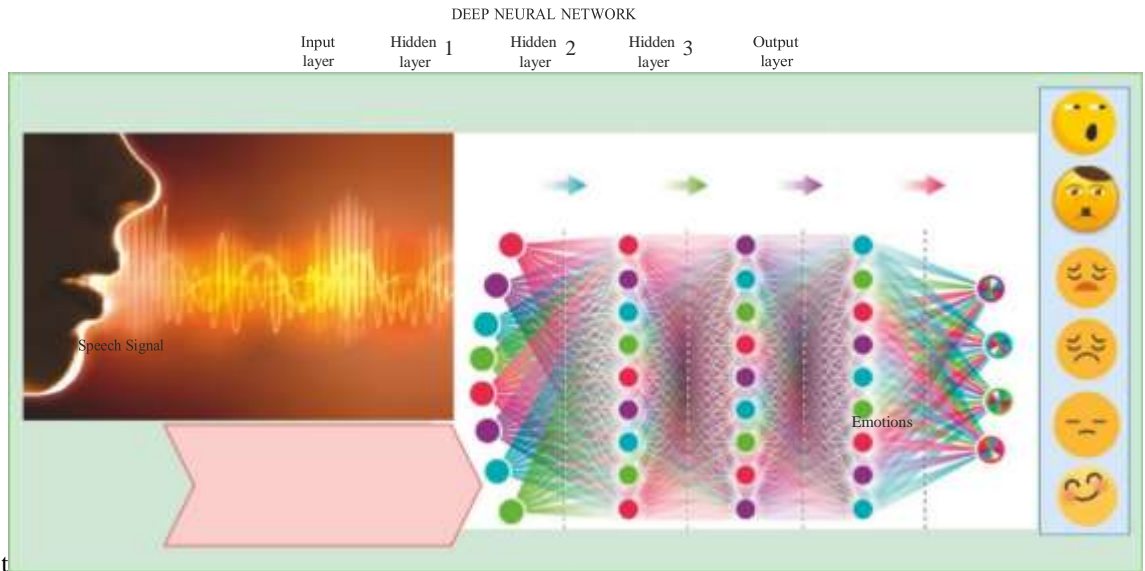
viable. Pattern analysis, audio recognition, computer vision, and image recognition are just some of the real-world applications that have reaped the benefits of deep learning techniques like convolutional neural networks, deep Boltzmann machines, deep belief networks, and stacked autoencoders. Recently, the aforementioned DL approaches suited to the FER problem were evaluated in depth by Li et al. [100]. As an introduction to CNN-based FER methods, Ginne et al. [5] have provided a useful summary. Extensive utilization of deep convolution has been shown in FER studies. Research into neural networks has mostly focused on becoming better at facial expression identification. The feasibility of a smaller CNN architecture with the same degree of accuracy should be carefully evaluated. Figure 8 shows how FER-based deep convolution neural networks may enhance deployment compatibility on memory-constrained devices, reduce costs, and increase dispersion by allowing for more effective dispersed training and a more adjustable parameter model.

shown by their use in a number of cutting-edge algorithms. Many FER competitions [101], including the previous year's EmotiW challenge, were won by a kind of CNN architecture with few layers. Facial emotion recognition has served the public well for decades prior to the field of deep learning breaking, and a group of brilliant researchers has tried to stay abreast of the current research efforts in that field, while others have undertaken to learn from its methods and discoveries. In recent times, many researchers offered novel and recurring practices for applying deep learning in order to security problems in an effort to enhance detection. Validation users currently do additional validation on a number of static or sequential databases before allowing their information to be used in a live database.

The VGG-16 model (developed by the University of Oxford's Visual Geometry Group (VGG)) may be considered a watershed moment in the history of deep CNN models [102]. It was pretrained using the ImageNet database

[103] to extract features from images that might be used to distinguish between image classes. Numerous recent studies show that VGG-16 performs well on image recognition and classification datasets from a variety of fields.

Marco et al. [104] proposed Deep Convolution Neural Networks (DCNNs) which are used in the cross-database search. After that, facial images had to be reduced to 48x48 pixels; the rest of the same pictures had to be searched for locations and landmarks to be extracted. Finally, they had augmented the database with additional data, and only then did they were able to create it. Subsequently, the data moves on to two classification stages where the softmax (SF) is expanded and fed into the fully connected softmax (XF) network after the first classification stage. To avoid over-fitting, they suggest using local CNNs in combination with convolutional layers that are fine-tuned for specific use cases. In [105], the authors have shown that the results prior to training were used to discover how to influence the final outcome. When it expanded, the first CNN



expansion, when it lowered the size to 32×32 and also used data

Figure 6: Speech emotion recognition [8].

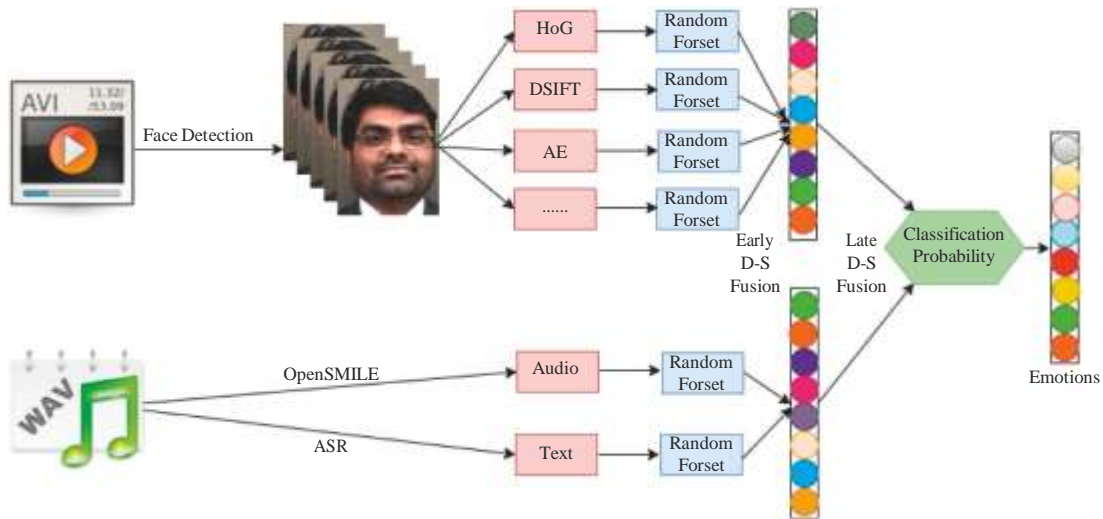


Figure 7: The framework for multimodal emotion recognition NN [94].

normalization with 8-connected pools followed by down-sampling (normalization of 32×32 to a 256 final dimension), and when that was done, cropping was employed. Gaining the most mass is something that happens only at the competition, so the athletes who have gained the most muscle will play in the games. The information used this third party search tool to assemble a total of three transparently accessible databases: the CK+ and JAFFE, as well as the BU3DF. One also discovers a wide range of beneficial practices when considering these studies, such as utilizing all of these techniques and products together; studies show the difference between these things yielding different results. The preprocessing techniques employed by Anil and his associates [8] have also been applied in the study by the authors. They are devising a new CNN face recognition algorithm for people who have not yet been recognized. They have two convolutions allowing layers and a dropout

layer that gives the net activation of one in order to predict more accurate results. They use maxing with one extra activation as well and a final convolution (expanding) layer in the last step to increase accuracy and flexibility. An important concern raised by Cai et al. [106] deals with the fact that the closing or the disappearance of small-town public libraries is managed by solving CNN, which employs Sparsity Batch Normalization. Dropout may be added to network building to help against overfitting and SBP (Support, Gradient, and Regularization) as a second stage to improve model generalization capacity, with the property of being used in networks twice (as a support for and then starting with 2-convol reg and ending with SGD), to strengthen the network. Li et al. [107] proposed using a CNN to tackle the facial distortion problem, in particular; the authors are doing so by first extracting the data from

the VGG network and then running the ACN. *Affect*. Also, architecture has been and has already been employed in the Affect Net, the RAF database, and the FED-RO database.

Yolcu and his colleagues [108] proposed that faces may be the most important aspects to focus on where Y could be realized in order to accurately take and record a single facial feature as small as an eyebrow or on top of the nose, like the face; the three microscopes have to examine a three-millimeter area. Before using the image for registration, they crop the face to avoid blurring. They work with only key-point facial regions until they are finished using the CNN to ensure that registration has been completed. Prior to this, the project being filmed in full-frame, the subjects' faces were exposed to be greatly improved and details increased; for example, expressions were added to them in order to show more complexity. There are more proven benefits; for example, studies show that utilizing photographs is a better approach for capturing the true appearance of your screen targets.

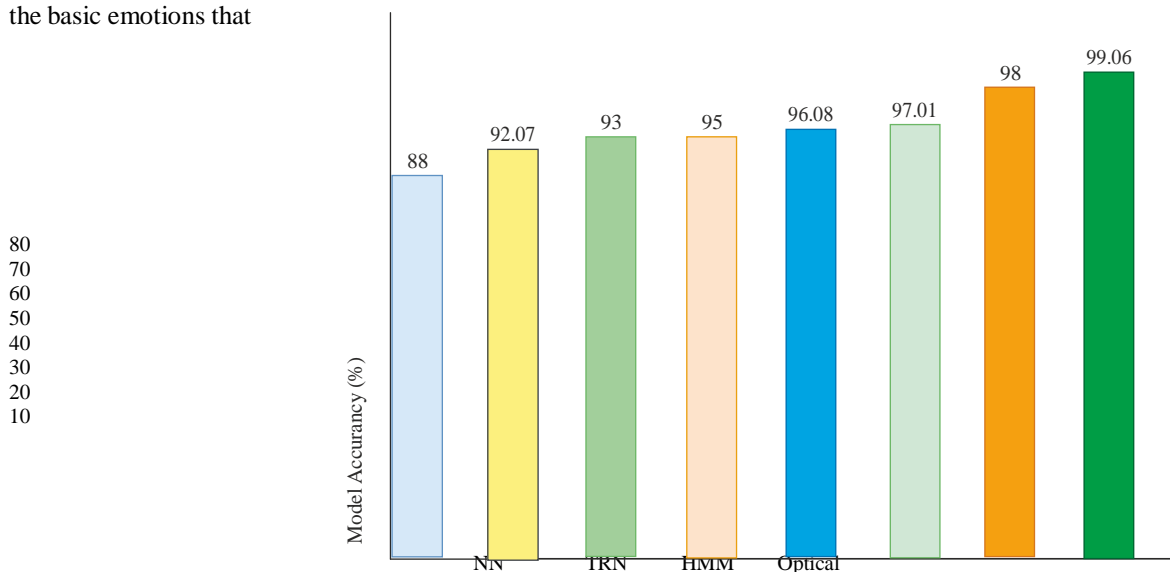
To figure out the significance of the CNN attributes in FER2013, researchers investigated and added to the already discovered findings of Agrawal et al. [109] (this also included research on Agrawalwarsh et al.) in 2019. Beyond that is an image memory pool at 64 64 pixel resolution, the network will have a certain type of an allowable number of

In the facial expression sequence level, DCNN was postulated by Liang [111] consisting of two deep layers, one of which handles spatial features and the other temporal features, which are treated as features that are then merged and expanded into vectors of 256 dimensions to form the large facial emotion category vector; that is, the expression differentiated into six basic emotions is utilized. They went through the Multitask Cascade Computational Net for face detection, after which they broadened the database with the technique of data augmentation. It is based on all of the scientists' opinions about classifying the basic emotions that

convolution layers, and ad hoc pooling will take the second position, followed by other admissible filters before classifiers. The results of the study demonstrate that determining models achieve a 61.23% and 63.77% of their accuracy using isolated units, compared to adjoining, or dropout models, but do not have well-connected layers.

New ideas were proposed by Deepak et al. [110], where they advise that residual blocks contain two channels, each of which has two consecutive convolution layers. These pretrained models after cropping and normalizing the images on the JAFP and CK + datasets before they go into training mode allow identifying and eliminating unwanted variations in intensity.

Kim et al. [118] compared the three emotional state models: they used CNN with LSTM to show facial expression variation in space and time. First, CNN expanded the facial state information into a spatial representation and used this representation to represent the temporal variation of facial information; afterward, CNN expanded temporal representations and preserved the spatial information. Furthermore, the authors [119] created a new network architecture known as Spatiotemporal Convolutional Network with Nested LSTM (STC-LSTM), which preserves both temporal and spatiotemporal features with a 3D-Cell T-LSTM style CNN.



DL Based FER Models

Figure 9: Deep learning-based FER models.

explain beforehand that joy, dread, surprise, sorrow, and apathy make up the emotional spectrum. Instead of sharing the findings and conclusions of other scholars, we showcase a variety of suggestions made by individuals in the know.

Ten. A Review of Related Literature and Discussion on FER In this study, we made it quite obvious that research on FER has to go beyond superficial learning methods. In its complete form, automated FER

consists of four primary data processing stages, many potential architectures, and, lastly, emotion recognition inside the core model. Overfitting the data and cropping and shrinking photos were only two of the many preprocessing strategies highlighted in this study to speed up the training and normalizing processes. Lopes and his coworkers believe that all of these methods have been explained well in their recent paper [120]. Models for FER based on deep learning are shown in Figure 9. Accuracy

was successfully accomplished by the numerous approaches and contributions discussed in this study. When it comes to neural network architectures, Moselhi et al. [121] demonstrated the critical importance of using neural networks and connection expansion layers. The authors Moham- madpour et al. [122] follow a long tradition of researchers who have chosen to extract AU from the face before attempting face-to-recognition. The purpose of this research is to learn more about the network and find out whether occlusion pictures exist. The use of the surplus blocks has been studied by Pise et al. [8]. While bigger eyes and smaller faces are all that can be shown in text imagery, the addition of an iconized face to the Networks, as shown by Yolu and Ayiv [108], enhance accuracy when working with low-resolution pictures. After carefully analyzing the recognition rate, we decided to add a two-concept CNN architecture extension by providing two additional feature articles to learn about the influence of CNN parameters. More than 90% of these experiments were successful in some way, and positive outcomes have been shown with most of the approaches used. Multiple

combinations were first presented by researchers focusing on spatial and temporal aspects; for example, CNN-L and 3D-CNN are often used to improve spatial features, but the combination also improves temporal features. According to the research done by Yu and coworkers, both the Kim et al. [118] and Liang et al. [111] approaches provide more accuracy than the one completed by the Kim group [118]. That's equivalent to an almost 99% increase in volume efficiency. Both temporal and spatial networks have proven effective in CNN applications. In order to get high accuracy in FER, the researchers opted for LSTM since it performs well not just with sequential data but also with time-dependent data. Softmax and Adam optimization are now the most challenging algorithms employed in CNN research and are used for parametric modeling. We also tested the model on several datasets to ensure consistency in findings and to verify the neural network design. With an eye on the architecture, database, and recognition rate addressed in the aforementioned papers, Table 1 summarizes the prior points.

TABLE 1: Comparison between FER models.

Approach	Technique	Groups	Sub	Authors	Acc (%)
DCBiLSTM	Fusion	6	123	Liang et al. [111]	99.6
Dist-based	Optical flow	5	8	Essa & pentland [112]	98
CNN	Facial AUs	7	123	Hashemi et al.,	97.01
SBN-CNN	Batch norm	7	10	Wei et al., [113]	96.8
Rule-based	Optical flow	6	32	Yacoob & davis [114]	95
HMM	2-D FT optical flow	6	4	Otsuka & Ohya	93
TRN	Relational reasoning	8	27	Pise et al. [8,115]	92.7
Rule-based	Parametric model	6	40	Black & Yacoob [116]	92
NN	Optical flow	2	32	Rosenblum et al. [117]	88

10. Conclusion;

Recent advances in FER were discussed, and the research given helped us to keep up with the field. Over the course of the last year or two, several scientists in and out of the lab have developed their own unique CNN architectures and reference datasets. The provision of both historical and experimental tables (spontaneous and lab) (emotion as reference) is necessary for reliable emotional identification. We also provide a debate that highlights the fact that robots can already detect more sophisticated emotions, suggesting the proliferation of human-machine cooperation. Future Plans An individual's emotional state can be gleaned via FER, but it can only learn the six fundamental emotions and neutral. It clashes with the reality of the situation, which includes more nuanced feelings. As a result, scientists will feel more motivated to increase the size of their datasets and design robust deep learning architectures that can identify both primary and secondary feelings. Furthermore, in the current period, emotion recognition has developed from a unimodal study into a multimodal analysis of a complex system. Multimodality, as shown by Leon et al. in [123], is essential for effective emotion recognition. Researchers are currently concentrating on creating and

commercializing effective multimodal deep learning architectures and databases, such the one that Zhang et al. [124] and Ringeval et al. [125] explored, which combines auditory and visual modalities with physiological data. Access to Information Contact the study's corresponding author if you'd want access to the data used to draw these conclusions. Interest Discrepancies There are no competing interests, as the authors have stated. Acknowledgments King Saud University, Riyadh, Saudi Arabia, is thanked for supporting this effort via the Researchers Sup- porting Project number (RSP-2022R426).

References

- "Automatic facial ex- pression recognition using facial animation parameters and multistream hmms," by P. S. Aleksic and A. K. Katsaggelos, was published in the January 2006 issue of IEEE Transactions on Information Forensics and Security.
- "Facial expression recognition using machine learning," by H. Chouhayebi, J. Riffi, M. A. Mahraz, Y. Ali, and T. Hamid, to appear in the Proceedings of the 2021 Fifth International Conference on Intelligent Computing in Data Sciences (ICDS), pp. 1-6, IEEE, Seoul, Korea, November 2021).
- (3) "Rela- tional reasoning using neural networks: a survey," International Journal of Uncertainty, Fuzziness, and Knowledge-

Based Systems, vol. 29, pp. 237-258, 2021, by Anil. Audumbar Pise, H. Vadapalli, and I. Sanders. "Not-So-CLEVR: learning same-different relations strains feedforward neural networks," by J. Kim, M. Ricci, and T. Serre, was published in 2018 in Interface focus, volume 8, issue 4. This is supported by the literature: [5] E. Sariyanidi, H. Gunes, and A. Cavallaro, "Automatic analysis of facial affect: a survey of registration, representation, and recognition," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 37, no. 6, pp. 1113-1133, 2014. Facial expression in affective disorders," What the Face Reveals: Basic and Applied Studies of Spontaneous Expression Using the Facial Action Coding System (FACS), vol. 2, pp. 331-342, 1997, E. Paul, D. Matsumoto, and V. F. Wallace. "Comprehensive database for facial expression analysis," by T. Kanade, J. F. Cohn, and Y. Tian, was published in the proceedings of the Fourth IEEE International Conference on Automatic Face and Gesture Recognition (Cat. No. PR00580), held in Grenoble, France, in March 2000 (see Customer satisfaction-aware operations planning," by J. Mei, K. Li, and K. Li. best multiserver setup for maximizing cloud computing profits," IEEE Transactions on Sustainable Computing, volume 2, issue 1, pages 17-29, 2017. University of California in San Francisco, "Emfac-7: emotional facial action coding system," vol. 2, no. 36, p. 1, 1983, V. F. Wallace, E. Paul, and others. xvi, 848, \$69.95] G. Schubert, "Human ethology and evolutionary episte-

mology: the strange case of dmEibesfeldthuman eieeyaPp. xvi, 848, \$69.95," Journal of Social and Biological Systems, volume 13, issue 4, pages 355-387, 1990. Based on the work of E. S. Jaha and L. Ghouti, "Color face recognition using quaternion pca," published in Proceedings of the 4th International Conference on Imaging for Crime Detection and Prevention 2011, pages 1-6, ICDP, 2011. Discussion of "Facial expressions of emotion: an old controversy and new findings" by Paul Ekman, Diana P. Perrett, and Harry D. Ellis, published in Philosophical Transactions of the Royal Society of London, Series A B, volume 335, page 69, 1992. [21] X. Zhou, K. Li, G. Xiao, Y. Zhou, and K. Li, "Top f probabilistic products queries," IEEE Transactions on Knowledge and Data Engineering, vol. 28, no. 10, pp. 2808-2821, 2016. For example, see [22] M. Bichsel and A. P. Pentland, "Human face recognition and the face image Set's topology," CVGIP: Image Understanding, vol. 59, no. 2, pages 254-261, 1994. IEEE Transactions on Parallel and Distributed Systems, volume 26, issue 11, pages 3040-3051, 2014. [23] K. Li, W. Ai, Z. Tang et al., "Hadoop recognition of bio- medical named entity using conditional random fields." According to [24] (X. Tang, K. Li, Z. Zeng, and B. Veeravalli), "A novel security- driven scheduling algorithm for precedence-constrained tasks in heterogeneous distributed systems," IEEE Transactions on Computers, vol. 60, no. 7, pages 1017-1029, 2010)