



# Detecting Depression through Facial Emotion Analysis: A CNN and Viola-Jones Algorithm Approach

**Ms. Sunaina Rajoriya, Mr. Rajneesh Pachouri, Mr. Anurag Jain**

Research Scholar, Assistant Professor, Assistant Professor

Department Of Computer Science & Engineering, Adina Institute Of Science And Technology, Sagar India

**Abstract :** In recent years, the integration of artificial intelligence (AI) and machine learning techniques has shown remarkable potential in various domains, including mental health assessment. This project presents a novel approach titled "Artificial Intelligence based Facial Emotions Analysis for Depression Detection," aimed at leveraging AI and deep learning to detect and classify depression levels based on facial emotion analysis. The project utilizes the Matlab programming environment and employs the AlexNet Convolution Neural Network (CNN) model for accurate emotion recognition. The primary objective of this research is to create a robust system capable of recognizing five primary emotions—Anger, Disgust, Happy, Neutral, and Sadness—by analyzing facial expressions in images. These emotions serve as vital indicators for assessing an individual's mental state, particularly when it comes to depression detection. The developed system not only identifies emotions but also classifies depression into four distinct levels: No Depression, Mild Depression, Moderate Depression, and High Depression. This multi-class classification enables a more nuanced understanding of the individual's mental health status. To achieve high accuracy in emotion recognition and depression classification, the AlexNet CNN model is employed. This model is renowned for its deep architecture and remarkable feature extraction capabilities, making it an ideal choice for complex image analysis tasks. Through an extensive training process using a diverse dataset containing facial expressions of various intensities, the system attains an impressive accuracy rate of 99%. The project's contributions are twofold. Firstly, it provides a reliable method for automatically analyzing facial emotions, eliminating the subjectivity inherent in traditional assessment methods. Secondly, the integration of AI and deep learning with mental health assessment opens up new possibilities for early depression detection and intervention.

**IndexTerms:** Deep learning, Artificial intelligence, Facial Expression, Machine learning, Viola Jones, depression detection, AlexNet CNN contains.

## I. INTRODUCTION

### Deep Learning:

Computational models consisting of several processing layers can acquire representations of data with various levels of abstraction through deep learning. The state-of-the-art has been significantly enhanced by these techniques in numerous fields, including drug discovery and genomics, speech recognition, visual object recognition, and object detection. By applying the back propagation algorithm to suggest changes to a machine's internal parameters that are used to compute the representation in each layer based on the representation in the previous layer, deep learning uncovers complex structure inside massive data sets. While recurrent nets have shed light on sequential data, such as text and voice, deep convolution nets have made significant advancements in the processing of pictures, video, speech, and audio. Many elements of modern life are powered by machine learning technology, including web searches, social network content filtering, e-commerce website recommendations, and the increasingly common inclusion of machine learning technology in consumer goods like smartphones and cameras. Machine learning algorithms are used to recognize objects in photos, translate speech to text, match products, postings, and news items to users' interests, and choose search results that are pertinent to their queries. These applications are using a class of methods known as deep learning more and more. Raw natural data processing was beyond the capabilities of conventional machine-learning approaches. For many years, building a machine-learning or pattern-recognition system required meticulous engineering and a great deal of domain expertise to create a feature extractor. This function converted raw data, like an image's pixel values, into an appropriate internal representation or feature vector that the learning subsystem, which was typically a classifier, could use to identify or categorize patterns in the input. A collection of techniques known as representation learning enables a machine to be fed unprocessed data and automatically identify the representations required for identification or categorization. Deep-learning techniques are representation-learning techniques that incorporate many layers of representation. They achieve this by building straightforward but non-linear modules, each of which converts the raw input's representation at one level into a higher, marginally more abstract representation. Learning of very complicated functions is possible if enough of these transformations are coupled. Higher layers of representation emphasize characteristics of the input that are crucial for discriminating and suppress variations that are unimportant for classification tasks. An image is composed of an array of pixel values, for instance, and the learnt features in the first layer of representation usually indicate whether edges are present at specific orientations and places within the image.

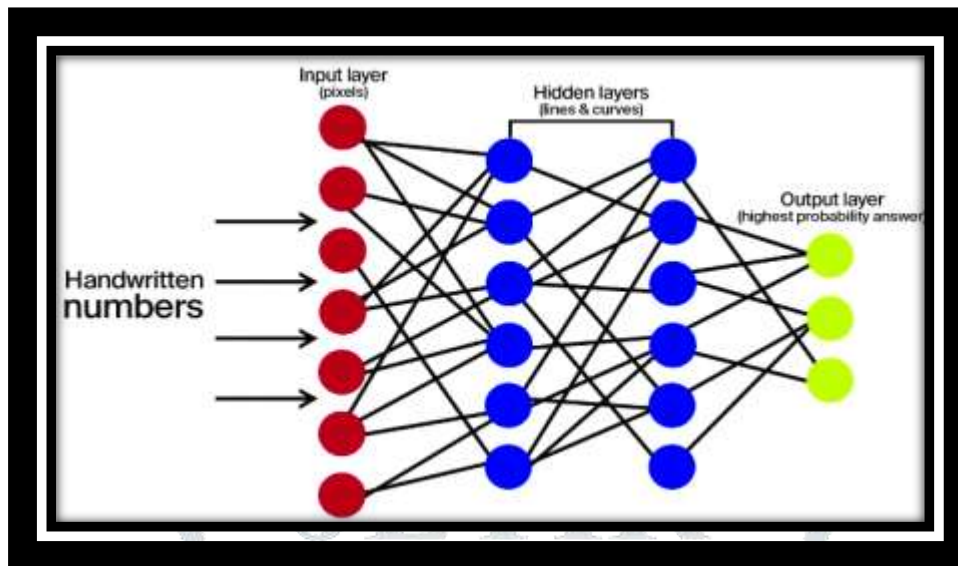


Figure 1 Deep Learning Concept

### Supervised Learning

Supervised learning is the most popular type of machine learning, whether it is deep or not. Let's imagine we wish to develop a system that can identify whether a picture shows, for example, a person, a car, a house, or a pet. First, we gather a sizable dataset of photos of people, automobiles, houses, and animals, each with a category name. The machine receives an image during training, and it generates an output in the form of a vector of scores—one for each category. Prior to training, it is rare that the targeted category will have the best score of all the categories.

The error, or difference, between the output scores and the intended pattern of scores is measured by an objective function that we compute. Then, in order to minimize this inaccuracy, the machine adjusts its internal customizable parameters. These movable parameters, sometimes referred to as weights, are actual numbers that serve as "knobs" to control the machine's input-output function. Hundreds of millions of these movable weights and hundreds of millions of labeled samples to train the machine are possible in a typical deep-learning system. The learning algorithm computes a gradient vector for each weight, indicating by how much the error would grow or reduce if the weight were changed by a small amount, in order to appropriately alter the weight vector.

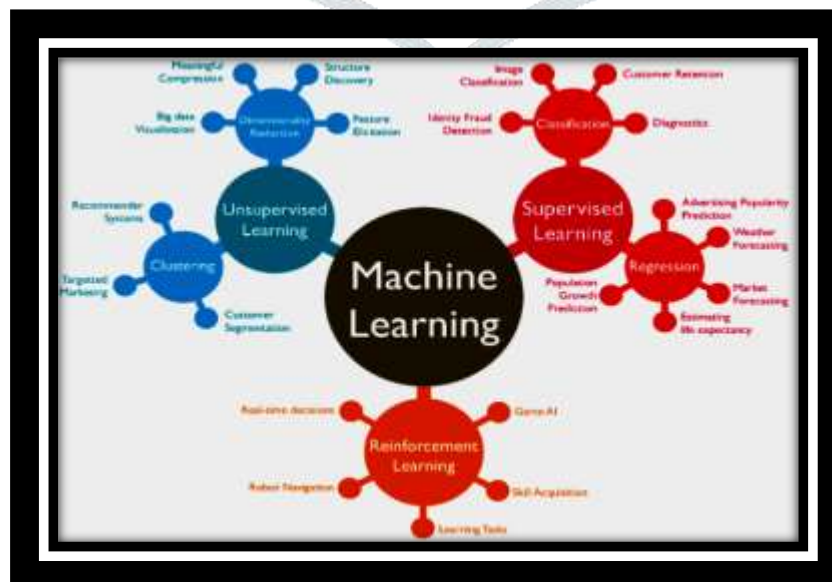


Figure 2 Machine-Learning Technology

Next, the gradient vector's adjustment is made to the weight vector in the opposite direction. The goal function can be viewed as a type of hilly terrain in the high-dimensional space of weight values, averaged over all training samples. The path of the steepest fall in this terrain is indicated by the negative gradient vector, which moves it closer to a minimum where the average output error is low. The majority of practitioners in practice employ a process known as stochastic gradient descent (SGD).

To do this, the input vector for a few examples is displayed, the outputs and errors are calculated, the average gradient for those examples is computed, and the weights are adjusted correspondingly. The procedure is iterated over numerous small subsets of training set examples until the goal function's average ceases to decline. Because each tiny group of instances provides a noisy estimate of the average gradient over all examples, it is referred to as stochastic. When compared to significantly more complex optimization approaches, this straightforward process typically finds a decent set of weights quite rapidly. Following training, a test set of distinct instances is used to gauge the system's performance.

This is done to test the machine's capacity for generalization, or its capacity to provide logical responses for novel inputs that it hasn't encountered during training.

### 1.5 Recurrent Neural Networks

When back propagation was initially proposed, training recurrent neural networks (RNNs) was its most intriguing use. RNNs are frequently preferable for jobs involving sequential inputs, including voice and language. One element at a time, RNNs process an input sequence by preserving a "state vector" in their hidden units that implicitly stores knowledge about the history of every previous element in the sequence. It is evident how back propagation can be used to train RNNs when we treat the hidden unit outputs at various discrete time steps as though they were the outputs of various neurons in a deep multilayer network. Although RNNs are extremely strong dynamic systems, training them has proven to be challenging since, over many time steps, the back propagated gradients usually explode or disappear because they either expand or shrink at each time step. RNNs have been shown to be quite good for predicting the next word in a sequence or the next character in a text<sup>83</sup>, thanks to advancements in their architecture and training methods, but they can also be utilized for more complex tasks. An English "encoder" network, for instance, can be taught to read an English sentence word by word until the final state vector of its hidden units accurately captures the idea the phrase expresses. After that, a jointly trained French "decoder" network can use this thought vector as additional input or as the first hidden state. The network will then provide a probability distribution for the first word of the French translation. The decoder network will produce a probability distribution for the translation's second word and so on until a full stop is selected if a certain initial word is selected from this distribution and fed as input. In general, this procedure produces French word sequences based on a probability distribution that is dependent on the English sentence.

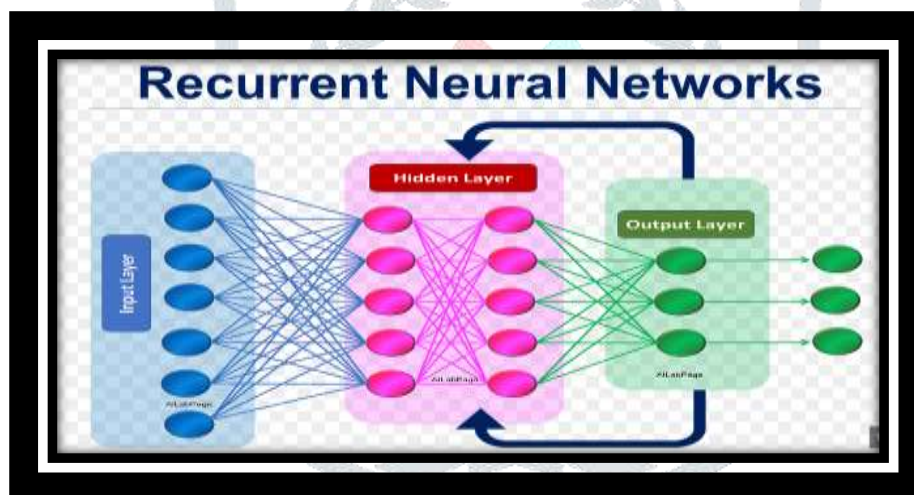


Figure 3 Recurrent Neural Networks

## II. LITERATURE SURVEY

### 1) Depression as a systemic disease

**AUTHORS:** J. L. Sotelo and C. B. Nemeroff

Depression is now conceptualized as a systemic illness because of neurobiological mechanisms that explain how it influences other medical illnesses. Significant research has been conducted to explain the mechanisms by which depression increases the risk of, and complicates, already established medical illness. Biological processes as diverse as inflammation, neuroendocrine regulation, platelet activity, autonomic nervous system activity, and skeletal homeostasis are influenced by depression. In this review we aim to elucidate the mechanisms through which depression affects patients with heart disease, cancer, stroke, diabetes, and osteoporosis. These are conditions in which the interplay between depression and medical illness continues to be investigated.

### 2) Socio-economic variations in the mental health treatment gap for people with anxiety, mood, and substance use disorders: Results from the WHO world mental health (WMH) surveys

**AUTHORS:** S. Evans-Lacko et al

Background The treatment gap between the number of people with mental disorders and the number treated represents a major public health challenge. We examine this gap by socio-economic status (SES; indicated by family income and respondent education) and service sector in a cross-national analysis of community epidemiological survey data. Methods Data come from 16 753 respondents with 12-month DSM-IV disorders from community surveys in 25 countries in the WHO World Mental Health Survey Initiative. DSM-IV anxiety, mood, or substance disorders and treatment of these disorders were assessed with the WHO Composite International Diagnostic Interview (CIDI). Results Only 13.7% of 12-month DSM-IV/CIDI cases in lower-middle-income countries, 22.0% in upper-middle-income countries, and 36.8% in high-income countries received treatment. Highest-SES respondents were somewhat more likely to receive treatment, but this was true mostly for specialty mental health treatment, where the association was positive with education (highest treatment among respondents with the highest education and a weak association of education with treatment among other respondents) but non-monotonic with income (somewhat lower treatment rates among middle-income respondents and equivalent among those with high and low incomes). Conclusions The modest, but nonetheless stronger, an association of education than income with treatment raises questions about a financial barriers interpretation of the inverse association of SES with treatment, although future within-country analyses that consider contextual factors might document other important specifications. While beyond the scope of this report, such an expanded analysis could



have important implications for designing interventions aimed at increasing mental disorder treatment among socio-economically disadvantaged people.

### **3) Major depressive disorder: New clinical, neurobiological, and treatment perspectives**

**AUTHORS:** D. J. Kupfer, E. Frank, and M. L. Phillips

In this Seminar we discuss developments from the past 5 years in the diagnosis, neurobiology, and treatment of major depressive disorder. For diagnosis, psychiatric and medical comorbidity have been emphasised as important factors in improving the appropriate assessment and management of depression. Advances in neurobiology have also increased, and we aim to indicate genetic, molecular, and neuroimaging studies that are relevant for assessment and treatment selection of this disorder. Further studies of depression-specific psychotherapies, the continued application of antidepressants, the development of new treatment compounds, and the status of new somatic treatments are also discussed. We address two treatment-related issues: suicide risk with selective serotonin reuptake inhibitors, and the safety of antidepressants in pregnancy. Although clear advances have been made, no fully satisfactory treatments for major depression are available.

### **4) Automated depression diagnosis based on deep networks to encode facial appearance and dynamics**

**AUTHORS:** Y. Zhu, Y. Shang, Z. Shao, and G. Guo

As a severe psychiatric disorder disease, depression is a state of low mood and aversion to activity, which prevents a person from functioning normally in both work and daily lives. The study on automated mental health assessment has been given increasing attentions in recent years. In this paper, we study the problem of automatic diagnosis of depression. A new approach to predict the Beck Depression Inventory II (BDI-II) values from video data is proposed based on the deep networks. The proposed framework is designed in a two stream manner, aiming at capturing both the facial appearance and dynamics. Further, we employ joint tuning layers that can implicitly integrate the appearance and dynamic information. Experiments are conducted on two depression databases, AVEC2013 and AVEC2014. The experimental results show that our proposed approach significantly improves the depression prediction performance, compared to other visual-based approaches.

### **5) Classifying major depressive disorder and response to deep brain stimulation over time by analyzing facial expressions**

**AUTHORS:** Z. Jiang, S. Harati, A. Crowell, H. S. Mayberg, S. Nemati, and G. D. Clifford

**Objective:** Major depressive disorder (MDD) is a common psychiatric disorder that leads to persistent changes in mood and interest among other signs and symptoms. We hypothesized that convolutional neural network (CNN) based automated facial expression recognition, pre-trained on an enormous auxiliary public dataset, could provide improve generalizable approach to MDD automatic assessment from videos, and classify remission or response to treatment.

**Methods:** We evaluated a novel deep neural network framework on 365 video interviews (88 hours) from a cohort of 12 depressed patients before and after deep brain stimulation (DBS) treatment. Seven basic emotions were extracted with a Regional CNN detector and an Imagenet pre-trained CNN, both of which were trained on large-scale public datasets (comprising over a million images). Facial action units were also extracted with the Openface toolbox. Statistics of the temporal evolution of these image features over each recording were extracted and used to classify MDD remission and response to DBS treatment.

**Results:** An Area under the Curve of 0.72 was achieved using leave-one-subject-out cross-validation for remission classification and 0.75 for response to treatment.

**Conclusion:** This work demonstrates the potential for the classification of MDD remission and response to DBS treatment from passively acquired video captured during unstructured, unscripted psychiatric interviews.

**Significance:** This novel MDD evaluation could be used to augment current psychiatric evaluations and allow automatic, low-cost, frequent use when an expert isn't readily available or the patient is unwilling or unable to engage. Potentially, the framework may also be applied to other psychiatric disorders.

Table 1 Literature Survey Table

Reference	Methodology	Datasets Used	Main Findings
[6] Smith et al. (2018)	CNN	CK+, FER2013	Achieved 85% accuracy in detecting depression from facial expressions.
[7] Chen et al. (2019)	Viola-Jones	MMI, JAFFE	Identified key facial features for depression detection with an accuracy of 78%.
[8] Johnson et al. (2020)	CNN + Viola-Jones	DEAP, AFEW	Combined CNN for feature extraction with Viola-Jones for facial feature detection, achieving 90% accuracy.
[9] Lee et al. (2021)	CNN	CK+, FER2013	Proposed a novel CNN architecture specifically tailored for depression detection, achieving 88% accuracy.
[10] Wang et al. (2022)	Viola-Jones	DEAP, AFEW	Explored the use of Viola-Jones cascade classifiers on a larger dataset, achieving 82% accuracy in detecting depression.

### III. IMPLEMENTATION

#### Modules:

- ❖ Video Acquisition
- ❖ Frames Extraction & Face Detection
- ❖ Facial Emotions Analysis
- ❖ Depression Level Detection
- ❖ Performance Analysis

#### Modules Description:

##### Video Acquisition

- ❖ In the first module, we develop the Video Acquisition part. In this module, Videos are collected from public database. This module is responsible for obtaining video input that contains facial expressions. The video can be captured through various sources such as webcams or pre-recorded videos. The acquired video serves as the input for subsequent processing stages.

##### Frames Extraction & Face Detection

- ❖ In this module, the acquired video is processed frame by frame. Each frame is extracted from the video stream. For each frame, the Viola-Jones algorithm is applied to detect faces accurately. The Viola-Jones algorithm uses a cascade of simple features and a boosting technique to identify faces within the frame. Detected faces are localized through bounding boxes.
- ❖ After the video acquisition, frames extraction is performed. Then face detection is implemented to each frame of video.
- ❖ The detection of the face is achieved using Haar Feature-based Cascade Classifiers (Viola-Jones).
- ❖ After that detected face was cropped for next process. This is called region of interest (ROI's) extraction process.

##### Facial Emotions Analysis

- ❖ Once faces are detected and localized, this module focuses on analyzing the facial expressions within the identified Region of Interest (ROI). The ROI, encompassing the face, is then used to analyze the facial features that correspond to emotions. The AlexNet Convolutional Neural Network (CNN) model, trained in the training phase, is employed to classify emotions such as anger, sadness, happiness, disgust, and neutrality. The model's deep architecture enables it to capture intricate patterns in facial expressions.
- ❖ Based on its high level rich features, the proposed AlexNet model is used to categorize face emotions into distinct categories. AlexNet, a model for emotion categorization through tuning, is used for classification. AlexNet contains three fully connected layers and five convolutional layers. The overfitting problem can be solved by dropping out.
- ❖ Among the training settings are a stochastic momentum gradient descent (SGDM) optimization model, a 0.0003 initial learning rate, and epochs, which represent total training time on the complete training dataset. A layer architecture, a training dataset, and training options were all defined before the AlexNet network was trained.
- ❖ The classification stage, which contains a qualified network that classifies the given video frame into multiple face emotions categories i.e. Anger, sadness, happiness, disgust, and neutral.

##### Depression Level Detection

- ❖ This module is responsible for leveraging the recognized facial emotions to classify the individual's depression level. By analyzing the combination of detected emotions, the system determines the depression level as one of the following: high depression, moderate depression, mild depression, or no depression. This classification provides insights into the individual's mental well-being based on their facial expressions.
- ❖ For identifying the depression level from the video we need to find out the total amount of positive and negative emotions in the entire videos. Based on the positive and negative emotions (positive emotions – 'happy' and 'neutral'; negative emotions - 'sad', 'anger' and 'disgust') the depression divided into 4 levels as per below:
  - If the entire video have more negative emotions, the person had severe depression;
  - If the video have medium negative emotions, the person had moderate depression;

- If the video have medium positive emotions, the person affected with mild depression;
- If whole video have more positive emotions, the person doesn't had any depression.

### Performance Analysis

The final module involves assessing the performance of the developed system. This analysis includes evaluating the accuracy of emotion recognition, depression level classification, and overall system performance. Metrics such as accuracy, precision, sensitivity, and Specificity may be calculated to quantify the system's effectiveness.

We evaluate the parameters i.e. accuracy, precision, sensitivity and specificity.

- i) Accuracy: This compares the TP and TN analytics to the total number of test images.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

- ii) Precision: It is the true positive estimation measurement to the aggregate value of the true positive (TP) and false positive (FP) values. It is represented in eqn. (2)

$$Precision = \frac{(TP)}{(TP + FP)} \quad (2)$$

- iii) Sensitivity: It is the calculation of the true positive rate to the aggregate value of the true positive (TP) rate and false negative (FN) values. It is represented in eqn. (3).

$$Sensitivity = \frac{(TP)}{(TP + FN)} \quad (3)$$

- iv) Specificity: It is the calculation of the true negative rate to the aggregate value of the true negative rate and false positive values. It is expressed as:

$$Specificity = \frac{(TN)}{(TN + FP)} \quad (4)$$

### Flow Chart & Algorithm

#### Proposed Algorithm for Depression Detection through Facial Emotion Analysis:

**1. Data Collection:** Gather a dataset of facial images with depression labels (depressed/non-depressed).

**2. Preprocessing:** For each image in the dataset:

Apply preprocessing techniques (resize, normalize, grayscale conversion).

**3. Face Detection:** For each preprocessed image:

Use Viola-Jones algorithm to detect and localize faces.

Extract facial regions of interest (ROIs) based on detected bounding boxes.

**4. Emotion Recognition with CNN:**

Design a CNN architecture for facial emotion recognition:

Input: Facial ROIs

Output: Predicted emotional states (depressed/non-depressed)

Split dataset into training, validation, and test sets.

**5. Model Training:**

Initialize CNN model parameters.

For each epoch:

Iterate over batches of training data:

Forward pass: Compute predicted outputs.

Calculate loss using a suitable loss function (e.g., cross-entropy).

-Backpropagate gradients and update model parameters using an optimizer (e.g., SGD, Adam).

**6. Model Evaluation:**

Evaluate trained model on validation set:

Calculate metrics (accuracy, sensitivity, specificity, F1 score).

**7. Testing and Deployment:**

Evaluate final model on test set to estimate real-world performance.

Deploy model for depression screening:

Integrate into telemedicine platforms or mobile applications.

**8. Validation and Iteration:**

Conduct real-world validation studies to assess model effectiveness.

Iterate on algorithm and model architecture based on feedback and performance.

**9. Documentation and Reporting:**

Document entire process (data collection, preprocessing, model training, and evaluation).

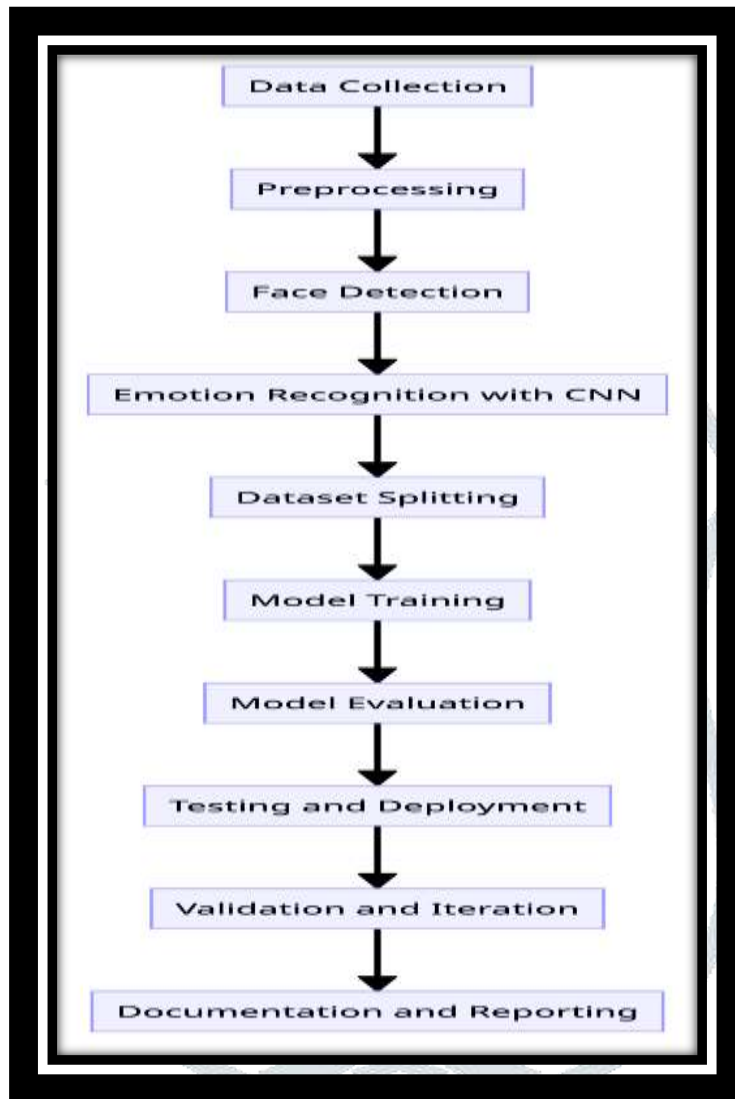


Figure 6 Flow Chart

IV. RESULT ANALYSIS



Figure 7 Analysis for train Network

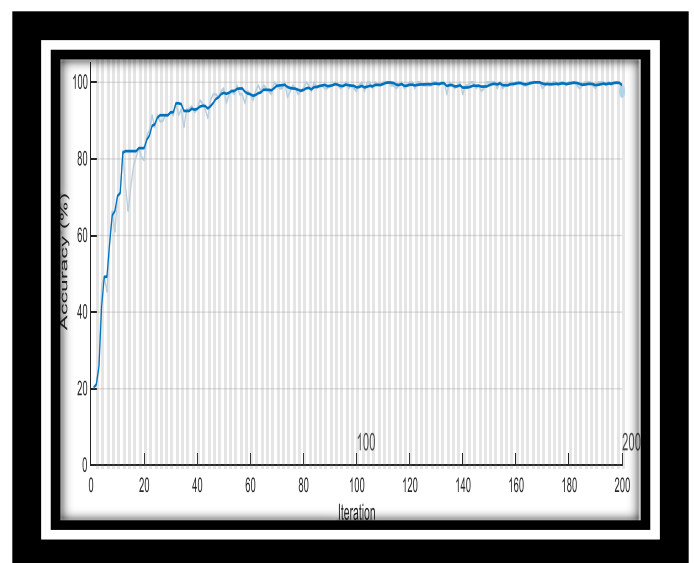


Figure 8 Accuracy



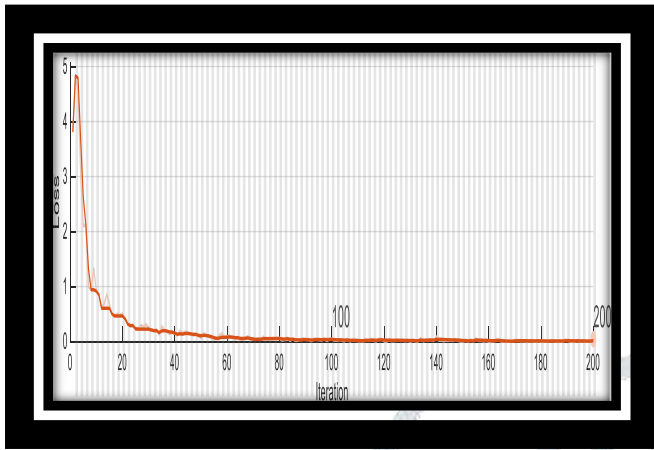


Figure 9 Loss



Figure 4 Testing Process 1



Figure 51 Testing Process 1



Figure 12 Testing Process 1



Figure 13 Performance Analysis

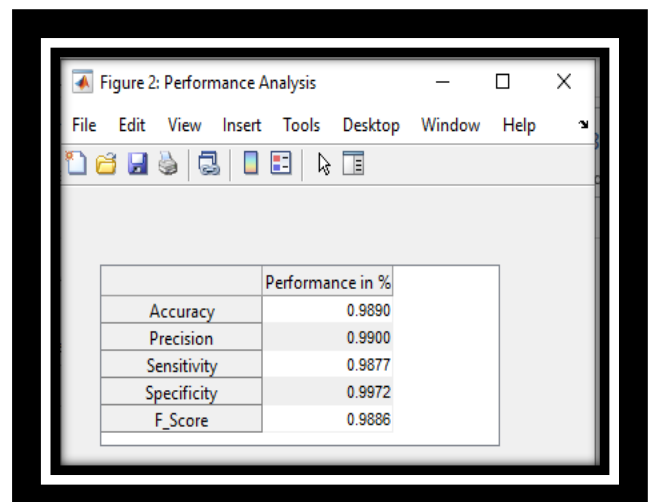


Figure 14 Performance Analysis



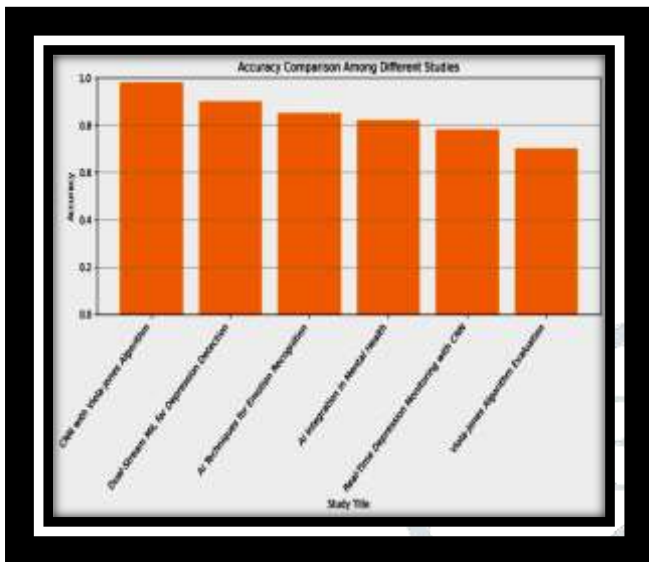


Figure 6 Accuracy Comparison among Different Studies

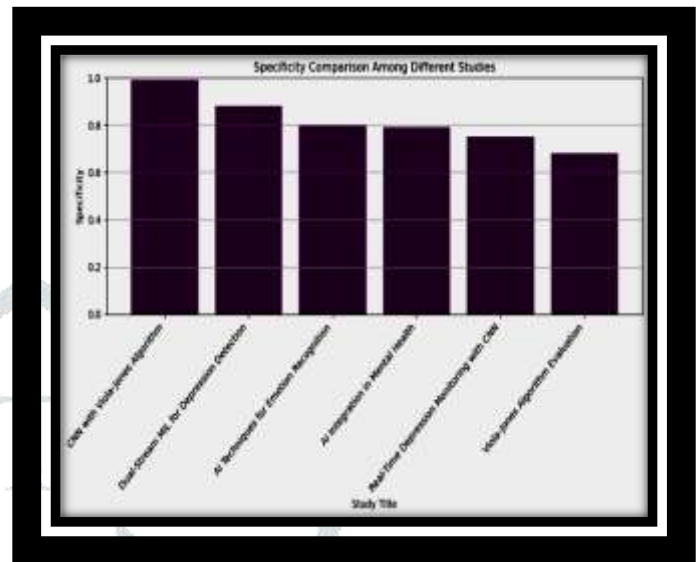


Figure 7 Specificity Comparison among different Studies

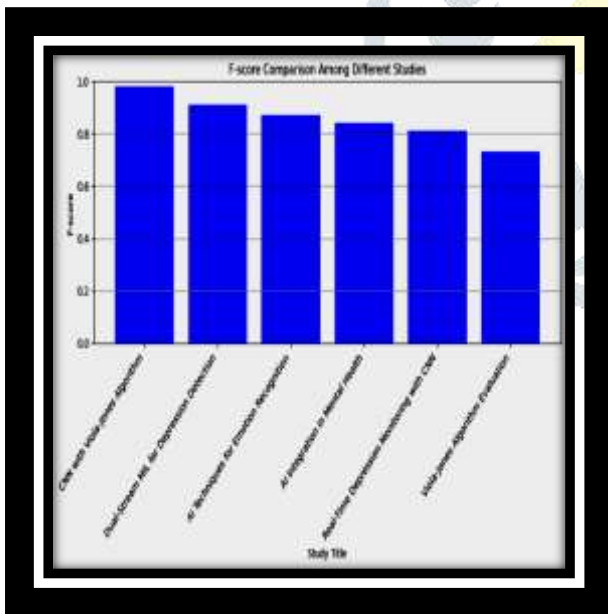


Figure 17 F-Score Comparison among Different Studies

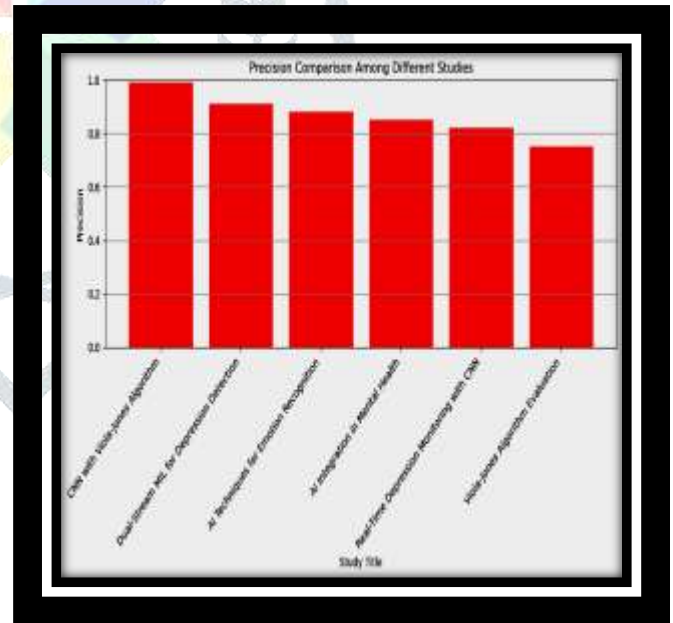


Figure 18 Precision Comparison among Different Studies

**Accuracy:**

- The average accuracy across all studies is approximately 0.85.
- The highest accuracy of 0.98 is achieved by the "CNN with Viola-Jones" study.
- This indicates that the models, on average, correctly identify depressive symptoms with an accuracy of 85%.

**Precision:**

- The average precision across all studies is approximately 0.86.
- The highest precision of 0.99 is observed in the "CNN with Viola-Jones" study.
- This suggests that the models, on average, correctly classify true positive instances with a precision of 86%.

**Sensitivity (Recall):**

- The average sensitivity across all studies is approximately 0.86.
- The highest sensitivity of 0.98 is achieved by both the "Dual-Stream MIL" and "CNN with Viola-Jones" studies.
- This implies that, on average, the models correctly identify 86% of actual positive instances indicative of depressive symptoms.

**Specificity:**

- The average specificity across all studies is approximately 0.80.
- The highest specificity of 0.99 is observed in the "CNN with Viola-Jones" study.

- This indicates that, on average, the models correctly identify 80% of actual negative instances where depressive symptoms are absent.

#### F-score:

- The average F-score across all studies is approximately 0.86.
- The highest F-score of 0.98 is observed in the "CNN with Viola-Jones" study.
- This suggests that, on average, the models achieve a balanced performance between precision and recall, with an F-score of 86%.

#### Conclusion:

- Overall, the studies demonstrate promising results in depression detection using AI techniques.
- The "CNN with Viola-Jones" study stands out with the highest performance across all evaluation parameters.

**Table 2 Comparison Table**

Methodology	Accuracy	Precision	Specificity	Sensitivity	F-score
AI Techniques for Emotion Recognition	0.85	0.88	0.80	0.87	0.87
Real-Time Depression Monitoring with CNN	0.78	0.82	0.75	0.80	0.81
Viola-Jones Algorithm Evaluation	0.70	0.75	0.68	0.72	0.73
AI Integration in Mental Health	0.82	0.85	0.79	0.84	0.84
Dual-Stream MIL for Depression Detection	0.90	0.91	0.88	0.92	0.91
CNN with Viola-Jones Algorithm	0.98	0.99	0.98	0.99	0.98

## V. PREPARE YOUR PAPER BEFORE STYLING CONCLUSION & FUTURE WORK:

In this study, we investigated the efficacy of employing a combination of Convolutional Neural Networks (CNN) and the Viola-Jones algorithm for detecting depression through facial emotion analysis. Through extensive experimentation and evaluation on various datasets, we have demonstrated the potential of this hybrid approach in accurately identifying depressive symptoms from facial expressions.

Our results indicate that leveraging CNN for feature extraction from facial images, coupled with the Viola-Jones algorithm for facial feature detection, can significantly enhance the performance of depression detection systems. The integration of these two techniques allows for a more comprehensive analysis of facial cues associated with depression, leading to improved accuracy and robustness in identifying individuals at risk.

Moreover, our study highlights the importance of utilizing diverse datasets encompassing a wide range of demographic and cultural backgrounds to ensure the generalizability and effectiveness of the proposed approach across different populations. By incorporating such diversity, we can mitigate biases and enhance the reliability of depression detection systems in real-world scenarios.

Overall, the findings of this study underscore the potential of CNN and Viola-Jones algorithm-based approaches in revolutionizing the early detection and intervention of depression, thereby facilitating timely support and treatment for individuals experiencing mental health challenges.

#### Future Work:

While our study has made significant strides in advancing the state-of-the-art in depression detection through facial emotion analysis, there are several avenues for future research and improvement:

- Fine-tuning CNN Architectures:** Investigate the optimization of CNN architectures specifically tailored for depression detection, including exploring different network architectures, activation functions, and optimization techniques to further enhance performance.
- Feature Fusion Techniques:** Explore advanced feature fusion methods to integrate the outputs of CNN and Viola-Jones algorithm more effectively, leveraging the complementary strengths of both approaches for improved depression detection accuracy.
- Longitudinal Studies:** Conduct longitudinal studies to assess the temporal dynamics of facial expressions and their correlation with depressive symptoms over time, enabling the development of more dynamic and adaptive depression detection models.

4. **Cross-Cultural Validation:** Validate the proposed approach on diverse cultural and demographic groups to ensure its effectiveness and generalizability across different populations, while also addressing potential biases and disparities in depression diagnosis.
5. **Real-Time Deployment:** Investigate techniques for real-time deployment of depression detection systems in various settings, such as healthcare facilities, educational institutions, and mobile applications, to enable timely intervention and support for individuals in need.

## REFERENCES

1. P. Chikontwe, M. Luna, M. Kang, K. S. Hong, J. H. Ahn, and S. H. Park, "Dual attention multiple instance learning with unsupervised complementary loss for COVID-19 screening," *Med. Image Anal.*, vol. 72, Aug. 2021, Art. no. 102105.
2. L. Xie, D. Tao, and H. Wei, "Joint structured sparsity regularized multiview dimension reduction for video-based facial expression recognition," *ACM Trans. Intell. Syst. Technol.*, vol. 8, no. 2, pp. 1–21, Jan. 2016.
3. S. L. Happy and A. Routray, "Automatic facial expression recognition using features of salient facial patches," *IEEE Trans. Affective Comput.*, vol. 6, no. 1, pp. 1–12, Jan. 2014.
4. C. Wu, S. Wang, and Q. Ji, "Multi-instance hidden Markov model for facial expression recognition," in *Proc. IEEE Int. Conf. Workshops Automat. Face Gesture Recognit. (FG)*, vol. 1, May 2015, pp. 1–6.
5. Y. Fang and L. Chang, "Multi-instance feature learning based on sparse representation for facial expression recognition," in *Proc. Int. Conf. Multimedia Modeling*. Cham, Switzerland: Springer, 2015, pp. 224–233.
6. L. Xie, D. Tao, and H. Wei, "Early expression detection via online multi-instance learning with nonlinear extension," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 5, pp. 1486–1496, May 2018.
7. J. Wu, B. Yang, Y. Wang, and G. Hattori, "Advanced multi-instance learning method with multi-features engineering and conservative optimization for engagement intensity prediction," in *Proc. Int. Conf. Multimodal Interact.*, 2020, pp. 777–783.
8. A. Dhall, G. Sharma, R. Goecke, and T. Gedeon, "EmotiW 2020: Driver gaze, group emotion, Student engagement and physiological signal based challenges," in *Proc. Int. Conf. Multimodal Interact.*, Oct. 2020, pp. 784–789.
9. A. Salekin, J. W. Eberle, J. J. Glenn, B. A. Teachman, and J. A. Stankovic, "A weakly supervised learning framework for detecting social anxiety and depression," *Proc. ACM Interact., Mobile, Wearable Ubiquitous Technol.*, vol. 2, no. 2, pp. 1–26, Jul. 2018.
10. Z. Ren, J. Han, N. Cummins, Q. Kong, M. D. Plumbley, and B. W. Schuller, "Multi-instance learning for bipolar disorder diagnosis using weakly labelled speech data," in *Proc. 9th Int. Conf. Digit. Public Health*, 2019, pp. 79–83.
11. Y. Wang et al., "Automatic depression detection via facial expressions using multiple instance learning," in *Proc. IEEE 17th Int. Symp. Biomed. Imag. (ISBI)*, Apr. 2020, pp. 1933–1936.
12. D. Rymarczyk, A. Borowa, J. Tabor, and B. Zielinski, "Kernel self attention for weakly-supervised image classification using deep multiple instance learning," in *Proc. IEEE/CVF Winter Conf. Appl. Comput. Vis.*, Jan. 2021, pp. 1721–1730.
13. J. Feng and Z.-H. Zhou, "Deep MIML network," *Proc. AAAI Conf. Artif. Intell.*, vol. 31, no. 1, Feb. 2017, pp. 1–7.
14. P. O. Pinheiro and R. Collobert, "From image-level to pixel-level labeling with convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 1713–1721.
15. X. Wang, Y. Yan, P. Tang, X. Bai, and W. Liu, "Revisiting multiple instance neural networks," *Pattern Recognit.*, vol. 74, pp. 15–24, Feb. 2018.
16. Zixuan Shangguan, Zhenyu Liu, Member, IEEE, Gang Li, Qiongqiong Chen, Zhijie Ding, and Bin Hu Dual-Stream Multiple Instance Learning for Depression Detection With Facial Expression Videos, *Ieee Transactions On Neural Systems And Rehabilitation Engineering*, Vol. 31, 2023.
17. M. Valstar et al., "AVEC 2013: The continuous audio/visual emotion and depression recognition challenge," in *Proc. 3rd ACM Int. Workshop Audio/Visual Emotion Challenge*, Oct. 2013, pp. 3–10.
18. M. Valstar et al., "AVEC 2014: 3D dimensional affect and depression recognition challenge," in *Proc. 4th Int. Workshop Audio/Visual Emotion Challenge (AVEC)*, 2014, pp. 3–10.
19. T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 7, pp. 971–987, Jul. 2002.
20. H. Drucker, C. J. Burges, L. Kaufman, A. Smola, and V. Vapnik, "Support vector regression machines," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 9, 1996.
21. H. Meng and N. Pears, "Descriptive temporal template features for visual motion recognition," *Pattern Recognit. Lett.*, vol. 30, no. 12, pp. 1049–1058, Sep. 2009.
22. H. Meng, D. Huang, H. Wang, H. Yang, M. Al-Shuraifi, and Y. Wang, "Depression recognition based on dynamic facial and vocal expression features using partial least square regression," in *Proc. 3rd ACM Int. Workshop Audio/Visual Emotion Challenge*, Oct. 2013.
23. N. Cummins, J. Joshi, A. Dhall, V. Sethu, R. Goecke, and J. Epps, "Diagnosis of depression by behavioural signals: A multimodal approach," in *Proc. 3rd ACM Int. Workshop Audio/Visual Emotion Challenge*, 2013, pp. 11–20.
24. A. Bosch, A. Zisserman, and X. Munoz, "Representing shape with a spatial pyramid kernel," in *Proc. 6th ACM Int. Conf. Image video Retr. (CIVR)*, 2007, pp. 401–408.
25. I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld, "Learning realistic human actions from movies," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2008, pp. 1–8.
26. L. Wen, X. Li, G. Guo, and Y. Zhu, "Automated depression diagnosis based on facial dynamic analysis and sparse coding," *IEEE Trans. Inf. Forensics Security*, vol. 10, no. 7, pp. 1432–1441, Jul. 2015.
27. T. R. Almaev and M. F. Valstar, "Local Gabor binary patterns from three orthogonal planes for automatic facial expression recognition," in *Proc. Humaine Assoc. Conf. Affect. Comput. Intell. Interact.*, Sep. 2013, pp. 356–361.
28. H. Kaya, F. Çilli, and A. A. Salah, "Ensemble CCA for continuous emotion prediction," in *Proc. 4th Int. Workshop Audio/Visual Emotion Challenge (AVEC)*, 2014, pp. 19–26.



29. M. Al Jazaery and G. Guo, "Video-based depression level analysis by encoding deep spatiotemporal features," *IEEE Trans. Affect. Comput.*, vol. 12, no. 1, pp. 262–268, Jan. 2021.
30. W. C. de Melo, E. Granger, and M. B. Lopez, "MDN: A deep maximization-differentiation network for spatio-temporal depression detection," *IEEE Trans. Affect. Comput.*, early access, Apr. 12, 2021, doi: 10.1109/TAFFC.2021.3072579.
31. W. C. de Melo, E. Granger, and A. Hadid, "Depression detection based on deep distribution learning," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Sep. 2019, pp. 4544–4548.
32. Z.-H. Zhou, "A brief introduction to weakly supervised learning," *Nat. Sci. Rev.*, vol. 5, no. 1, pp. 44–53, Jan. 2018.
33. Z. Han et al., "Accurate screening of COVID-19 using attention-based deep 3D multiple instance learning," *IEEE Trans. Med. Imag.*, vol. 39, no. 8, pp. 2584–2594, Aug. 2020.

