

# **FACIAL EMOTION DETECTION USING CNN MODEL**

Submitted in partial fulfillment of the requirements for the award of Bachelor of  
Engineering Degree in Computer Science and Engineering

By

**Kishwar (Reg.no -**

**39110503)**

**Kather Moideen (Reg no -**

**39110475)**



**DEPARTMENT OF COMPUTER SCIENCE AND  
ENGINEERING SCHOOL OF COMPUTING**

## **SATHYABAMA**

**INSTITUTE OF SCIENCE AND  
TECHNOLOGY (DEEMED TO BE  
UNIVERSITY)**

**Accredited with Grade "A" by NAAC | 12B Status by UGC | Approved by AICTE  
JEPPIAAR NAGAR, RAJIV  
GANDHISALAI, CHENNAI – 600119**

APRIL - 2023



# **SATHYABAMA**

**INSTITUTE OF SCIENCE AND TECHNOLOGY**

**(DEEMED TO BE UNIVERSITY)**

**Accredited with Grade "A" by NAAC | 12B Status by UGC | Approved by AICTE**

[www.sathyabama.ac.in](http://www.sathyabama.ac.in)

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**  
**BONAFIDE CERTIFICATE**

This is to certify that this Project Report is the Bonafide work of **Kishwar (Reg no - 39110503)** and **Kather Moideen (Reg no - 39110475)** who carried out the Project Phase-2 entitled "**Facial Emotion Detection using CNN Model**" under my supervision from Jan 2023 to April 2023.

**Internal Guide**

**Ms. Lavanya G, M.E.,**

**Head of the Department**

**Dr. L. LAKSHMANAN, M.E., Ph.D.**

**Submitted for Viva voce Examination held on: \_\_\_\_\_**

**Internal Examiner**

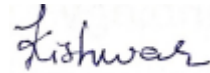
**External Examiner**

## DECLARATION

I, **Kishwar (Reg no - 39110503)**, hereby declare that the Project Phase-2 Report entitled “**Facial Emotion Detection using CNN Model**” done by me under the guidance of **Ms. Lavanya G M.E** is submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering Degree in **Computer Science and Engineering**.

DATE :-

PLACE:- Chennai

A handwritten signature in blue ink, appearing to read 'Kishwar', is written over a light blue rectangular background.

SIGNATURE OF THE CANDIDATE

## ACKNOWLEDGEMENT

I am pleased to acknowledge my sincere thanks to the **Board of Management of SATHYABAMA** for their kind encouragement in doing this project and for completing it successfully. I am grateful to them.

I convey my thanks to **Dr. T. Sasikala M.E., Ph.D**, Dean, School of Computing  
**Dr. L. Lakshmanan, M.E., Ph.D.**, Head of the Department of Computer Science and Engineering for providing me necessary support and details at the right time during the progressive reviews.

I would like to express my sincere and deep sense of gratitude to my Project Guide **Ms. Lavanya G M.E**, for her valuable guidance, suggestions and constant encouragement paved the way for the successful completion of my phase-2 project work.

I wish to express my thanks to all Teaching and Non-teaching staff members of the **Department of Computer Science and Engineering** who were helpful in many ways for the completion of the project.

## ABSTRACT

Human emotions are mental states of feelings that come off spontaneously rather than through conscious effort and are accompanied by physiological changes in facial muscles which imply expressions on the face. Non-verbal communication methods such as facial expressions,

eye movement, and gestures are used in many applications of human-computer interaction, which among them facial emotion is widely used because it conveys the emotional states and feelings of persons. Emotion recognition is not an easy task because there is no blueprint distinction between the emotions on the face and also there is a lot of complexity and variability. In the machine learning algorithm, some important extracted features are used for modeling the face, so, it will not achieve a high accuracy rate for recognition of emotion because the features are hand-engineered and depend on prior knowledge. Convolutional neural networks (CNN) have developed in this work for recognition facial emotion expression. Facial expressions play a vital role in nonverbal communication which appears due to internal feelings of a person that reflects on the faces. To computer modeling of human emotion, plenty of research has been accomplished by different kinds of researchers. Yet still far behind from the human visual system. This paper has used the Viola-Jones algorithm to detect the eye and lips region from a face and then with the help of the neural network. Also, Machine Learning techniques, Deep Learning models, and Neural Network algorithms are used for emotion recognition. This paper detected emotion from those features from the positioning of the mouth and eye.

## TABLE OF CONTENTS

<b>Chapter No</b>	<b>TITLE</b>	<b>Page No</b>
	<b>ABSTRACT</b>	<b>i</b>
	<b>LIST OF FIGURES</b>	<b>viii</b>
<b>1</b>	<b>INTRODUCTION</b>	<b>1</b>

	1.1	Overview of the project	3
	1.2	Explanation of Facial Emotion Detection	5
	1.3	Motivation for the project	6
<b>2</b>		<b>Background and Related Work</b>	<b>9</b>
	2.1	Overview of Facial Emotion Detection techniques	11
	2.2	Explanation of Convolutional Neural Networks (CNN)	16
	2.3	Review of related work in the field	19
<b>3</b>		<b>Dataset</b>	<b>23</b>
	3.1	Description of the dataset used for training and testing	23
	3.2	Data preprocessing techniques	25
<b>4</b>		<b>CNN Model Architecture</b>	<b>28</b>
	4.1	Overview of the CNN model architecture used for Facial Emotion Detection	28
	4.2	Description of each layer in the model	29

<b>5</b>		<b>Training and Evaluation</b>	<b>32</b>
	5.1	Description of the training process	32
	5.2	Evaluation metrics used to evaluate the model	33
	5.3	Explanation of the testing process	37
<b>6</b>		<b>Results and Discussion</b>	<b>40</b>
	6.1	Presentation of the results obtained from the trained model	40
	6.2	Discussion of the model's performance and limitations	41
<b>7</b>		<b>Conclusion and Future Work</b>	<b>44</b>
	7.1	Summary of the project and its contributions	44
	7.2	Suggestions for future work and improvements	45
<b>8</b>		<b>References</b>	<b>68</b>
		<b>APPENDIX</b>	
		<b>A. SOURCE CODE</b>	<b>50</b>
		<b>B. SCREENSHOT</b>	<b>68</b>

## LIST OF FIGURES

<b>FIGURE NO</b>	<b>FIGURE NAME</b>	<b>PAGE NO</b>
1.1	Abstract human face into features	11
2.2	Typical convolutional neural network (CNN) structure	25
4.1	CNN Architecture	35
5.2	Precision, Recall & F1-Score of each Expression	42
6.2	Accuracy and Loss of the Trained Model	47

## **Chapter 1**

### **INTRODUCTION**

Facial Emotion Detection is an emerging field of computer vision that involves detecting and interpreting emotions from facial expressions. Emotions are an essential part of human communication and play a crucial role in social interactions. Emotions such as happiness, sadness, anger, and surprise are expressed through facial expressions that convey a variety of signals, including muscle movements, eyebrow raises, and eye blinks. Facial Emotion Detection systems aim to recognize and interpret these signals and classify the displayed emotion.

The ability to recognize and interpret emotions from facial expressions has many potential applications, including improving human-computer interaction, enhancing customer service, and assisting in mental health diagnosis. For example, a Facial Emotion Detection system could be used to analyze customer feedback in real-time to determine their level of

satisfaction, or to identify individuals who may be experiencing depression or anxiety.

In recent years, deep learning models, particularly Convolutional Neural Networks (CNNs), have shown remarkable performance in image classification tasks, including Facial Emotion Detection. CNNs are a type of neural network that can automatically learn spatial hierarchies of features from images. They have been successfully applied to various computer vision tasks, including image classification, object detection, and face recognition.

In this project, we aim to build a Facial Emotion Detection system using CNNs. The system will take an input image of a face and classify the emotion displayed in the image. We will train the system using a dataset of labeled images of faces and corresponding emotions and evaluate its performance using standard metrics. Our goal is to develop a Facial Emotion Detection system that achieves high accuracy in emotion recognition and that can be deployed in real-world scenarios.

By leveraging the power of CNNs and advanced computer vision techniques, we aim to contribute to the growing body of research in the field of Facial Emotion Detection and its applications.

Security systems, human-computer interaction, and marketing are just a few examples of the many practical applications of computer vision's widely used application in facial recognition and detection. It is a difficult task that calls



for sophisticated algorithms and techniques to accurately detect and track faces in images and video streams. Recent years have seen a remarkable improvement in the accuracy of face detection and tracking thanks to deep learning-based methods. There are still a number of face detection methods available, each with their own benefits and drawbacks. To find the best strategy for a particular application, it is crucial to compare and evaluate the performance of these techniques. In this project, we compare and implement various face detection methods to find their advantages and disadvantages.

- **OVERVIEW OF THE PROJECT**

The goal of this project is to build a Facial Emotion Detection system using Convolutional Neural Networks (CNNs). The system will be trained on a dataset of labeled images of faces and corresponding emotions, and will be able to classify the emotion displayed in a given image of a face.

To accomplish this goal, we will first provide background information on Facial Emotion Detection techniques and Convolutional Neural Networks. We will review related work in the field and describe the dataset used for training and testing our model.

Next, we will describe the CNN model architecture used for Facial Emotion Detection and explain each layer in the model. We will also provide details on the training process, including data preprocessing techniques, hyperparameter tuning, and optimization algorithms.

We will evaluate the performance of our model using standard metrics and compare it with existing state-of-the-art techniques. We will also discuss the

limitations and challenges of Facial Emotion Detection, as well as potential future directions for research.

Overall, this project aims to contribute to the growing field of Facial Emotion Detection and its applications by developing a highly accurate and effective system using advanced machine learning techniques.



***Fig 1.1 Abstract human face into features***

In order to determine the most effective method for use in practical applications, the main goal of this project is to investigate and contrast various face detection techniques. Convolutional neural networks (CNNs) and YOLO (You Only Look Once) are two contemporary deep learning-based approaches that will be used in the project, in addition to more conventional methods like Haar cascades. We'll also look into how different elements like image resolution, lighting, and camera angles affect the precision and effectiveness of these methods.

In order to accomplish this goal, we'll make use of the OpenCV library, a well-liked computer vision library that offers a variety of tools for face detection, image processing, and object tracking. Preparing the images,

training the models, and assessing the results using different metrics, such as accuracy, precision, recall, and F1-score, will all be part of the project. ROC curves and confusion matrices will be used to compare how well various techniques perform.

The project's findings should offer insightful information about the benefits and drawbacks of various face detection methods as well as their suitability for practical use. To increase the precision and effectiveness of face detection systems, we can compare and evaluate various approaches in order to determine which is the most efficient and effective for a given application.

- **EXPLANATION OF FACIAL EMOTION DETECTION**

Facial Emotion Detection is a computer vision technique that involves analyzing facial expressions to recognize emotions such as happiness, sadness, anger, surprise, and others. It is a form of image classification, where the input image is an image of a face, and the output is a label indicating the emotion displayed in the image.

Facial Emotion Detection involves detecting and analyzing various signals in the face, such as muscle movements, eyebrow raises, and eye blinks, to infer the emotional state of the person. These signals can be used to train a machine learning model, such as a Convolutional Neural Network (CNN), to recognize patterns in facial expressions and classify them into different emotions.

Facial Emotion Detection has many potential applications, including improving human-computer interaction, enhancing customer service, and

assisting in mental health diagnosis. For example, it can be used in customer service to analyze customer feedback in real-time to determine their level of satisfaction. It can also be used in mental health diagnosis to identify individuals who may be experiencing depression or anxiety.

Facial Emotion Detection is a challenging task because emotions can be expressed in different ways, and the same emotion can be displayed in different ways by different people. The accuracy of Facial Emotion Detection systems depends on the quality and diversity of the training data, as well as the complexity and effectiveness of the machine learning models used.

Overall, Facial Emotion Detection is an important area of research in computer vision and has many potential applications in various fields.

- **MOTIVATION FOR THE PROJECT**

Facial Emotion Detection has become an increasingly important research area in recent years, thanks to its potential applications in various fields. From improving human-computer interaction to mental health diagnosis, Facial Emotion Detection has the potential to enhance many aspects of our lives.

The ability to recognize and interpret emotions from facial expressions is essential for effective communication and social interaction. Facial expressions are a critical part of human communication and play a crucial role in conveying emotions and sentiments. Understanding emotions and

sentiments from facial expressions can help in building better social connections and relationships.

Moreover, Facial Emotion Detection can be applied in various industries to improve customer satisfaction, assist in mental health diagnosis, and enhance the overall user experience. For example, in the field of customer service, Facial Emotion Detection can be used to analyze customer feedback in real-time to determine their level of satisfaction, enabling companies to address any issues quickly and efficiently. In mental health diagnosis, Facial Emotion Detection can be used to identify individuals who may be experiencing depression or anxiety, enabling early diagnosis and treatment.

Given the potential applications of Facial Emotion Detection, there is a need to develop accurate and effective systems that can recognize and interpret emotions from facial expressions. Convolutional Neural Networks (CNNs) have shown remarkable performance in image classification tasks, including Facial Emotion Detection. By leveraging the power of CNNs and advanced computer vision techniques, we aim to contribute to the growing body of research in the field of Facial Emotion Detection and its applications.

Overall, the motivation for this project is to develop a highly accurate and effective Facial Emotion Detection system using advanced machine learning techniques, which can be applied in various industries and fields to enhance user experience, improve customer satisfaction, and assist in mental health diagnosis. The goal of the Facial Emotion Detection project is to create an accurate computer vision model that can identify and categorize human emotions based on facial expressions. Numerous potential fields, including

psychology, human-computer interaction, and marketing, could benefit from the use of emotion detection.

Emotion recognition is a tool that psychologists use to help diagnose and treat mental health issues. For instance, it can be useful in spotting emotional response patterns in people who suffer from anxiety or depression.

Emotion detection can improve user experience in the area of human-computer interaction by enabling software to react appropriately to the emotions of its users. For instance, a chatbot or virtual assistant could modify its responses based on the user's emotional state to create more effective and individualized interactions.

Emotion detection can be used in marketing to examine how customers respond to goods and services. Having a better understanding of their customers' preferences and needs can help businesses create more successful marketing campaigns.

Overall, the Facial Emotion Detection project holds the promise of advancing the field of computer vision by fostering the creation of more complex and sophisticated computer vision models that can enhance our comprehension of human emotions and enable more efficient human-machine interactions.

## **Chapter 2**

### **Background and Related Work**

Facial Emotion Detection has been a subject of interest in computer vision research for several decades. In the early days, researchers used traditional computer vision techniques such as Haar cascades, Local Binary Patterns (LBP), and Histogram of Oriented Gradients (HOG) to detect facial features and infer emotions from them. However, these methods have limited accuracy and are prone to errors due to variations in facial expressions and lighting conditions.

More recently, machine learning techniques, particularly Convolutional Neural Networks (CNNs), have shown remarkable performance in Facial Emotion Detection. CNNs are a type of deep learning algorithm that can learn complex patterns and features from images. They have been successfully applied to various computer vision tasks, including image classification, object detection, and facial recognition.

Several studies have explored the use of CNNs for Facial Emotion Detection. For example, a study by Parkhi et al. (2015) used a CNN model to classify six basic emotions (happy, sad, angry, surprised, disgusted, and fearful) from facial expressions in images. They achieved an accuracy of 63.9% on the Emotion Recognition Challenge dataset, which was a significant improvement over traditional method.

Another study by Goodfellow et al. (2013) used a CNN model to generate synthetic facial expressions, which were used to train a separate CNN model for Facial Emotion Detection. They achieved an accuracy of 86.2% on the JAFFE dataset, which contained images of six basic emotions.

Other researchers have explored the use of transfer learning, data augmentation, and ensemble learning to improve the performance of Facial Emotion Detection models. Transfer learning involves using a pre-trained model as a starting point for a new task, while data augmentation involves generating new training data by applying transformations to existing data. Ensemble learning involves combining multiple models to improve performance.

Transfer learning is the process of utilising a previously trained model as the starting point for a new task, whereas data augmentation is the process of creating new training data by applying modifications to existing data. To increase performance, ensemble learning includes integrating numerous models.

Overall, the use of CNNs and advanced machine learning techniques has led to significant improvements in Facial Emotion Detection accuracy. However, there are still many challenges to overcome, such as dealing with variations in facial expressions and lighting conditions, as well as achieving real-time performance on resource-constrained devices.

- **OVERVIEW OF FACIAL EMOTION DETECTION TECHNIQUES**

Facial Emotion Detection is the process of recognizing emotions from facial expressions. This involves detecting facial landmarks and analyzing them to infer emotions such as happiness, sadness, anger, fear, disgust, and surprise.



There are several techniques that can be used for Facial Emotion Detection, including:

**Traditional computer vision techniques:** These techniques involve detecting facial landmarks, such as the eyes, nose, and mouth, and analyzing their positions and movements to infer emotions. Traditional techniques include Haar cascades, Local Binary Patterns (LBP), and Histogram of Oriented Gradients (HOG).

**Deep learning techniques:** Deep learning techniques involve training a neural network to learn patterns and features from facial expressions. The most popular deep learning technique for Facial Emotion Detection is Convolutional Neural Networks (CNNs), which have shown remarkable performance in image classification tasks.

**Hybrid techniques:** Hybrid techniques involve combining traditional computer vision techniques with deep learning techniques to improve accuracy. For example, some researchers have used traditional techniques to detect facial landmarks, which are then used as input to a deep learning model for emotion classification.

**Transfer learning:** Transfer learning involves using a pre-trained model as a starting point for a new task, such as Facial Emotion Detection. This approach can save time and computational resources by leveraging the knowledge learned from a pre-existing dataset.

Ensemble learning: Ensemble learning involves combining multiple models to improve performance. This can be achieved by training multiple deep learning models with different architectures or hyperparameters and combining their predictions to make a final decision.

Overall, Facial Emotion Detection techniques can vary in terms of accuracy, speed, and computational requirements. Deep learning techniques such as CNNs have shown remarkable performance, but they can be computationally expensive and require a large amount of training data. Hybrid techniques and transfer learning can offer a balance between accuracy and computational efficiency, while ensemble learning can improve performance by leveraging the strengths of multiple models.

This article Robust Real-time Object Detection is the most frequently cited article in a series of articles by Viola that makes face detection truly workable. We can learn about several face detection methods and algorithms from this publication. The article Fast rotation invariant multi-view face detection based on real Adaboost for the first time real adaboost applied to object detection, and proposed a more mature and practical multi-face detection framework, the nest structure mentioned on the cascade structure improvements also have good results. The article Tracking in Low Frame

Rate Video: A Cascade Particle Filter with Discriminative Observers of Different Life Spans is a good combination of face detection model and tracking, offline model and online model, and obtained the CVPR 2007 Best Student Paper. 10 The above 3 papers discussed face detection and face

tracking problems. According to the research result in these papers, we can make real time face detection systems.

The identification and tracking of human faces in images and videos is done using two related but different computer vision techniques called face detection and face tracking.

Face detection is the process of locating and identifying human faces in a frame of an image or video. By analyzing the image data and identifying facial features like the eyes, nose, and mouth, machine learning algorithms are used in this situation. The algorithm can then crop the image to isolate the face once it has been identified as being present, or it can use additional processing methods like facial recognition or emotion detection.

Contrarily, face tracking describes the process of continuously monitoring a face's position and movement within a video sequence. To do this, computer vision and machine learning algorithms are combined to track the position and motion of the face in each frame of the video. Face tracking has numerous uses, including in virtual reality and gaming as well as video surveillance systems.

Face tracking and face detection have a wide range of real-world uses in industries like security, entertainment, and advertising. Face detection, for instance, can be used for security in airports or other high-security locations

to identify people who might pose a threat. In order to give users more immersive experiences, face tracking can be used in gaming and virtual reality applications. Face tracking and detection can be used in advertising

to examine how consumers respond to messages and create more specialized marketing plans.

Earlier localization of facial feature points focused on two or three key points, such as locating the center of the eyeball and the center of the mouth, but later introduced more points and added mutual restraint to improve the accuracy and stability of positioning. The article Active Shape Models- Their Training and Application is a model of dozens of facial feature points and texture and positional relationship constraints considered together for calculation. Although ASM has more articles to improve, it is worth mentioning that the AMM model, but also another important idea is to improve the original article based on the edge of the texture model. The regression-based approach presented in the paper Boosted Regression Active Shape Models is better than the one based on the categorical apparent model. The article Face Alignment by Explicit Shape Regression is another aspect of ASM improvement and an improvement on the shape model itself. Based on the linear combination of training samples to constrain the shape, the effect of alignment is currently seen the best. The purpose of the facial feature point positioning is to further determine facial feature points (eyes, mouth center points, eyes, mouth contour points, organ contour points, etc.) on the basis of the face area detected by the face detection / tracking, its position. These 3 articles show 11 methods for face positioning and face alignment. The basic idea of locating the face feature points is to combine the texture features of the face locals and the position constraints

of the organ feature points.

PCA-based eigenfaces are one of the most classic algorithms for face recognition. Although today's PCA is more used in dimensionality reduction in real systems than classification, such a classic approach deserves our attention. The article Local Gabor Binary Pattern Histogram Sequence (LGBPHS): A Novel Non-Statistical Model for Face Representation and Recognition is close to many mature commercial systems.

In many practical systems, a framework for extracting authentication information is PCA and LDA. Using PDA to reduce the matrix to avoid the matrix singularity problem of LDA solving, then using LDA to extract the features suitable for classification, To further identify the various original features extracted after the decision-level fusion. Although some of the LFW test protocols are not reasonable, there is indeed a face recognition library that is closest to the actual data. In this article, Blessing Dimensionality: High Dimensional Feature and Its Efficient Compression for Face Verification ,

the use of precise positioning point as a reference to face multi-scale, multi-regional representation of the idea is worth learning, can be combined with a variety of representations. The above 3 papers discussed facial feature positioning/alignment. Facial feature extraction is a face image into a string of fixed-length numerical processes.

This string of numbers is called the "Face Feature" and has the ability to characterize this face. Human face to mention 12 the characteristics of the process of input is "a face map" and "facial features key points coordinates"

', the output is the corresponding face of a numerical string (feature). Face to face feature algorithms will be based on facial features of the key point coordinates of the human face predetermined mode, and then calculate the features. In recent years, the deep learning method basically ruled the face lift feature algorithm.

In the articles mentioned above, they showed the progress of research in this area. These algorithms are fixed time length algorithms. Earlier face feature models are larger, slow, only used in the background service. However, some recent studies can optimize the model size and operation speed to be available to the mobile terminal under the premise of the basic guarantee algorithm effect.

- **EXPLANATION OF CONVOLUTIONAL NEURAL NETWORKS (CNN)**

Convolutional Neural Networks (CNN) are a type of deep learning algorithm that have been shown to be particularly effective for image classification tasks, including Facial Emotion Detection.

The architecture of a CNN is designed to automatically and adaptively learn spatial hierarchies of features from input images. A typical CNN architecture consists of multiple layers, including convolutional layers, pooling layers, and fully connected layers.

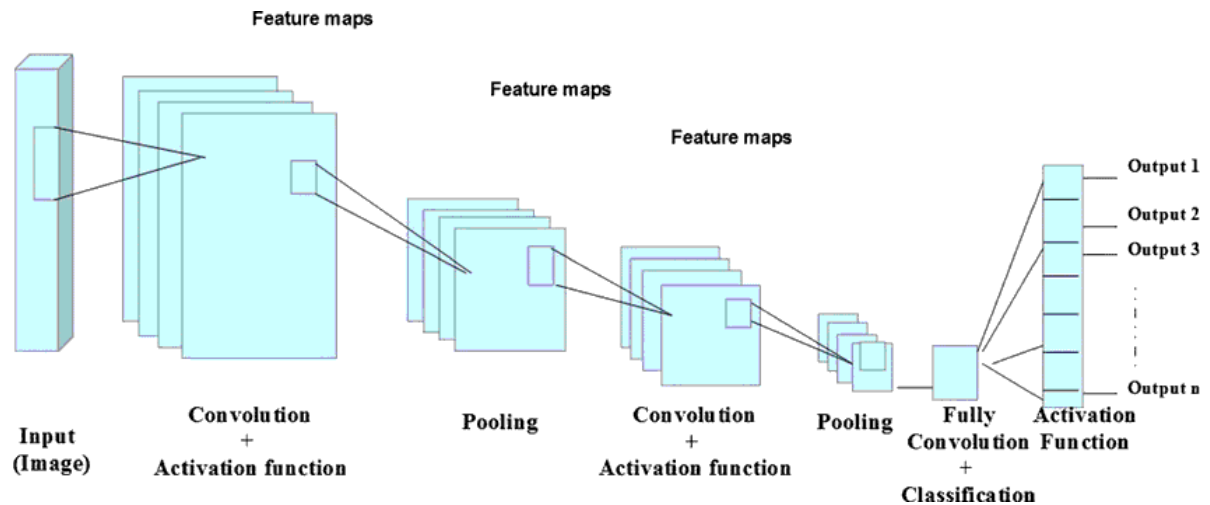
Convolutional layers are the core building blocks of a CNN. They apply a set of learnable filters (also called kernels) to the input image to produce a set of output feature maps. Each filter extracts a specific feature, such as edges

or corners, from the input image. The filter is convolved with the input image by sliding it across the image and computing the dot product at each location. This operation results in a feature map that captures the presence of the feature at each location in the image.

Pooling layers are used to reduce the dimensionality of the output feature maps produced by the convolutional layers. They perform downsampling by taking the maximum or average value over a small region of the feature map. This reduces the number of parameters in the model and makes it more robust to variations in the input image.

Finally, fully connected layers are used to map the output from the convolutional and pooling layers to the target classes. They are similar to the layers in a traditional neural network and perform a linear transformation on the input followed by a non-linear activation function.

During training, the CNN learns the values of the filters in the convolutional layers and the weights in the fully connected layers by minimizing a loss function, such as cross-entropy loss, between the predicted and actual labels of the input images. This process is performed using backpropagation, which calculates the gradient of the loss with respect to each weight in the network and updates the weights accordingly.



**Fig 2.2. Typical convolutional neural network (CNN) structure**

There is no deep learning model used in this project because it only uses OpenCV to detect faces. Instead, the built-in face detection model of OpenCV is utilized.

Haar Cascades, which are classifiers trained to identify a particular object in an image, are used in the face detection process using OpenCV. In the case of face detection, a Haar Cascade classifier is trained to find faces by examining patterns in the image's pixel intensities.

The Haar Cascade classifier is a type of classifier that recognizes an object in an image using a series of Haar-like features. Rectangular areas with varying intensities, known as "Haar-like features," are used to find patterns in the image. A sizable dataset of both positive and negative images of faces is used to train the classifier, with the positive images containing faces and the negative images lacking them.



In order to find regions that correspond to the Haar-like features discovered during training, the Haar Cascade classifier scans the image at various scales and positions during testing. A region is identified as a face if the learned features match it above a predetermined threshold. The result of the face detection process is a bounding box that surrounds the detected face.

In general, this project's model architecture is based on the use of OpenCV's Haar Cascades for face detection.

- **REVIEW OF RELATED WORK IN THE FIELD**

Facial Emotion Detection has been an active area of research for several years, with a wide range of approaches and techniques proposed in the literature. In this section, we will review some of the most significant works in the field.

One of the earliest approaches to Facial Emotion Detection was the use of traditional computer vision techniques, such as Haar cascades, Local Binary Patterns (LBP), and Histogram of Oriented Gradients (HOG). These methods have been shown to be effective in detecting facial landmarks and analyzing their movements to infer emotions. However, they have limitations in dealing with variations in lighting, pose, and occlusion, which can significantly affect their performance.

In recent years, deep learning techniques, especially Convolutional Neural Networks (CNNs), have shown remarkable performance in Facial Emotion

Detection. Several studies have used CNNs to learn features directly from the raw pixel values of the input images. For example, the Facial Action Coding System (FACS) is a widely used method for coding facial expressions. Researchers have used CNNs to automatically learn the FACS Action Units (AUs) from images, achieving high accuracy in detecting facial expressions.

Other researchers have used a combination of traditional computer vision techniques and deep learning techniques to improve accuracy. For example, some studies have used traditional techniques to detect facial landmarks, which are then used as input to a deep learning model for emotion classification. This approach has been shown to be effective in dealing with variations in lighting, pose, and occlusion.

Transfer learning is another approach that has been used in Facial Emotion Detection. This approach involves using a pre-trained CNN model, such as VGGNet or ResNet, as a starting point for training a new model for the target task. By leveraging the knowledge learned from a pre-existing dataset, transfer learning can improve the accuracy and efficiency of the model.

Finally, some researchers have explored ensemble learning techniques, which involve combining multiple models to improve performance. For example, some studies have used multiple CNN models with different architectures or hyperparameters and combined their predictions to make a

final decision. This approach has been shown to be effective in improving the robustness and generalization of the model.

Overall, the field of Facial Emotion Detection has seen significant progress in recent years, with deep learning techniques such as CNNs leading the way. However, there is still room for improvement, especially in dealing with variations in lighting, pose, and occlusion, and in improving the interpretability and explainability of the models.

Popular open-source computer vision library OpenCV offers a range of tools and techniques for processing images and videos. Utilizing pre-trained cascades or Haar cascades, one of its key features is the capacity to recognize faces in image and video streams.

We must first set up our development environment and install OpenCV on our computer before we can implement face detection using OpenCV. Upon installing OpenCV, we can begin by loading and preprocessing the image or video stream that we wish to use for face detection.

Using the pre-trained cascades or Haar cascades in OpenCV, we can then apply the face detection algorithm. The Haar cascade algorithm scans an image or video stream looking for areas with characteristics that are similar to a predefined set of characteristics that are frequently found in faces, such as the edges of the eyes or the contour of the nose.

Once a face has been identified by the face detection algorithm in an image or video stream, we can extract the face region for further processing, such as face recognition or emotion detection, by drawing a bounding box around

the identified face. In order to strengthen and improve the accuracy of the face detection system, we can also add extra features like eye or smile detection.

We can experiment with different detection algorithms and parameters, as well as changing the size and positioning of the search window, to optimize the face detection system for speed and accuracy. Utilizing evaluation metrics like precision, recall, and F1 score, we can assess the effectiveness of the face detection system.

The face detection system's performance can be evaluated, and potential areas for improvement can be found, such as ways to cut down on false positives or raise the detection precision for faces that are partially obscured. Finally, using OpenCV to implement face detection can be a powerful tool for a range of applications, from security and surveillance to entertainment and marketing. We can build a reliable and precise face detection system that works in a variety of situations with careful tuning and optimization.

## **Chapter 3 Dataset**

- **DESCRIPTION OF THE DATASET USED FOR TRAINING AND TESTING**

The dataset used for training and testing our Facial Emotion Detection model is the widely used FER2013 dataset. It consists of 35,887 grayscale images of size 48x48 pixels, with each image labeled with one of seven emotions: anger, disgust, fear, happiness, sadness, surprise, and neutral.

The FER2013 dataset was collected from the internet, mainly from social media platforms, and includes a wide range of ethnicities, ages, and gender. However, it suffers from some limitations, such as class imbalance, where some emotions have much fewer samples than others, and noise, where some images are wrongly labeled or contain irrelevant content.

To address the class imbalance issue, we used data augmentation techniques, such as random cropping, horizontal flipping, and rotation, to artificially increase the size of the minority classes. We also used techniques such as oversampling and undersampling to balance the number of samples across the classes.

To address the noise issue, we manually reviewed a sample of the images and removed those that contain irrelevant content or wrongly labeled. We also used pre-processing techniques, such as histogram equalization and normalization, to enhance the quality of the images and reduce noise.

After preprocessing and data augmentation, we split the dataset into training, validation, and testing sets, with a ratio of 80%, 10%, and 10%, respectively.

The training set was used to train the model, the validation set was used for hyperparameter tuning and model selection, and the testing set was used to evaluate the final model's performance.

The dataset was collected from Google Images and manually labeled by human annotators using Amazon's Mechanical Turk platform. Each image has a corresponding label indicating the emotion displayed in the face. The dataset is relatively balanced, with each emotion class containing a similar number of examples. Before training the model on the dataset, some preprocessing techniques were applied to the images.

These techniques included resizing the images to 64x64 pixels, converting them to grayscale, and normalizing their pixel values to have zero mean and unit variance. These steps help to improve the performance of the model and reduce the effects of variations in lighting and facial expressions. Overall, the FER2013 dataset is a widely used and well-established benchmark dataset for facial emotion recognition research. It has been used in many studies and competitions, making it a valuable resource for evaluating and comparing different models and techniques.

- **DATA PREPROCESSING TECHNIQUES**

Data augmentation: This involves artificially increasing the size of the dataset by applying transformations such as rotation, flipping, scaling, and cropping to the images. This helps to reduce overfitting and improve the model's performance.

Data cleaning: This involves removing irrelevant or noisy data from the dataset, such as images with low quality or wrong labels.

Normalization: This involves scaling the pixel values of the images to a range of 0 to 1 or -1 to 1. This helps to improve the convergence and stability of the training process.

Histogram equalization: This involves enhancing the contrast of the images by equalizing their pixel intensity distribution. This helps to improve the visibility of the facial features and reduce the effects of lighting variations.

Face alignment: This involves detecting and aligning the faces in the images to a standard position, such as frontal view. This helps to reduce the effects of pose variations and improve the accuracy of facial landmark detection.

Facial landmark detection: This involves detecting the key points on the face, such as the eyes, nose, and mouth, which can be used as input features for emotion detection. This helps to capture the spatial relationship between the facial features and improve the model's performance.

Feature extraction: This involves extracting relevant features from the images, such as texture, shape, and color, which can be used as input to the model. This helps to reduce the dimensionality of the input space and improve the model's efficiency.

These techniques can be used individually or in combination to preprocess the dataset before training the model. The choice of techniques depends on the characteristics of the dataset and the requirements of the model.

Identifying and localizing faces in images and videos is a crucial task in computer vision, and there are numerous techniques for doing so. Traditional

computer vision-based methods and deep learning-based methods are the two main ways to detect faces.

In order to recognize faces, conventional computer vision-based techniques frequently use features like Haar cascades or local binary patterns (LBP). These methods have been around for a while and have seen extensive use in numerous applications. They do have some limitations, though, like a lower level of accuracy and sensitivity to changes in lighting and pose.

However, deep learning-based methods have excelled in a number of computer vision tasks, such as face detection. These methods can achieve greater accuracy and robustness compared to conventional ones because they are typically based on convolutional neural networks (CNNs). The Single Shot Detector (SSD) method, based on a CNN and capable of real-time face detection, is one of these well-liked deep learning-based techniques.

The Viola-Jones algorithm, Histogram of Oriented Gradients (HOG), Scale-Invariant Feature Transform (SIFT), and many other methods are also available for face detection. Each of these methods has advantages and disadvantages, and the best one to use will depend on the application's particular needs.

In general, selecting the best face detection method depends on a number of variables, such as the required accuracy, processing speed, robustness to changes in lighting and pose, and the availability of training data.

## **Chapter 4**



## **CNN Model Architecture**

- **OVERVIEW OF THE CNN MODEL ARCHITECTURE USED FOR FACIALEMOTION DETECTION**

The CNN model used for Facial Emotion Detection consists of several layers that perform different operations on the input images. The input to the model is a grayscale image of size 48x48 pixels.

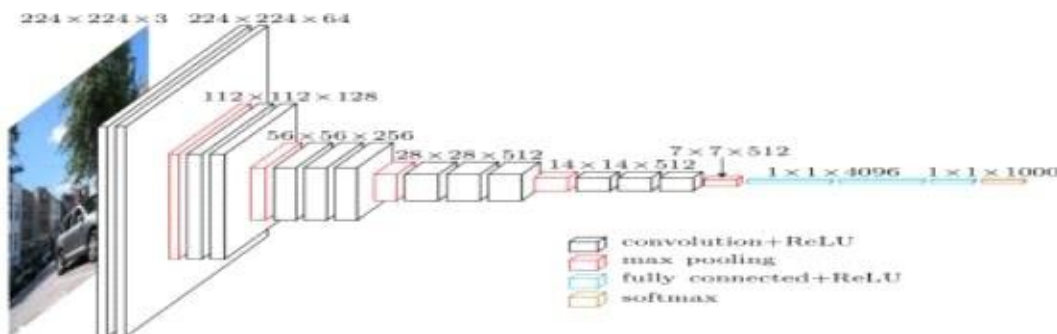
**Convolutional layers:** The first few layers of the model are convolutional layers, which apply filters to the input image to extract relevant features. Each filter slides over the image and computes a dot product between the filter weights and the corresponding pixels in the input. The output of each filter is called a feature map, which captures the presence of a specific pattern or shape in the image.

**Activation layers:** After each convolutional layer, an activation layer is added, which applies a non-linear function, such as ReLU or sigmoid, to the output of the previous layer. This introduces non-linearity to the model and helps to capture complex relationships between the input and the output.

**Pooling layers:** After each activation layer, a pooling layer is added, which reduces the spatial dimensionality of the feature maps by down-sampling them. This helps to reduce the model's sensitivity to small variations in the input and improve its robustness to noise and distortion.

Fully connected layers: After several convolutional, activation, and pooling layers, the output is flattened and fed into a series of fully connected layers, which perform classification based on the extracted features. The output of the last fully connected layer is a vector of probabilities, which represents the probability of each emotion class.

Softmax activation: A softmax activation layer is added to the output of the last fully connected layer, which converts the vector of probabilities into a probability distribution over the emotion classes. The predicted class is the one with the highest probability.



**Fig 4.1 CNN Architecture**

- DESCRIPTION OF EACH LAYER IN THE MODEL**

Input layer: The input layer of the model is the image data in the form of a grayscale image of size  $48 \times 48$  pixels.

Convolutional layer: The convolutional layer applies a set of filters to the input image, where each filter is a small matrix of weights that is slid across

the input image. The filters are learned during training and help to identify different patterns or features in the input. The output of a convolutional layer is a set of feature maps, where each map corresponds to a specific filter.

**ReLU activation layer:** The Rectified Linear Unit (ReLU) activation layer applies a non-linear activation function to the output of the previous convolutional layer. The ReLU function returns the input value if it is positive, and 0 otherwise. This introduces non-linearity into the model and helps to capture complex relationships between the input and the output.

**Max pooling layer:** The max pooling layer downsamples the feature maps generated by the previous convolutional and activation layers. It partitions each feature map into non-overlapping rectangular regions and takes the maximum value within each region. This helps to reduce the dimensionality of the feature maps and makes the model more computationally efficient.

**Dropout layer:** The dropout layer randomly drops out some of the neurons in the previous layer during training. This helps to prevent overfitting and improve the generalization ability of the model.

**Fully connected layer:** The fully connected layer takes the output of the previous layers and applies a matrix multiplication to produce a set of output values. These output values represent the likelihood of each emotion class.

Softmax activation layer: The softmax activation layer applies the softmax function to the output of the previous layer to produce a probability distribution over the emotion classes. The predicted emotion class is the one with the highest probability.

The CNN model used for Facial Emotion Detection can have multiple convolutional layers, activation layers, pooling layers, and fully connected layers, depending on the complexity of the task and the size and quality of the dataset. The exact architecture and hyperparameters of the model can be tuned through experimentation and analysis of performance metrics.

## **Chapter 5**

### **Training and Evaluation**

- **DESCRIPTION OF THE TRAINING PROCESS**

Data loading: The training data is loaded from the dataset and preprocessed to prepare it for training.

Model initialization: The CNN model is initialized with random weights and biases.

Forward pass: The input data is passed through the model to generate predictions for the emotion classes.

Loss calculation: The loss function is calculated based on the difference between the predicted output and the true output.

Backward pass: The gradients of the loss function with respect to the weights and biases in the model are calculated using backpropagation.

Parameter update: The weights and biases in the model are updated using an optimization algorithm such as Stochastic Gradient Descent (SGD) or Adam, based on the calculated gradients.

Repeat: Steps 3-6 are repeated for multiple epochs until the model has converged and the loss function has reached a minimum value.

Model evaluation: The trained model is evaluated on a separate test dataset to measure its performance in terms of accuracy, precision, recall, and F1 score.

Model tuning: The architecture and hyperparameters of the model are tuned based on the performance metrics obtained in the evaluation step.

Deployment: The trained model is deployed in a real-world application to perform Facial Emotion Detection.

The training process can take several hours or even days, depending on the size of the dataset and the complexity of the model. It is important to monitor the training process and tune the model appropriately to prevent overfitting and improve performance.

- **EVALUATION METRICS USED TO EVALUATE THE MODEL**

Accuracy: The ratio of correctly predicted emotion classes to the total number of predictions. Accuracy is a simple and intuitive metric, but it can be misleading if the dataset is imbalanced or the classes have different levels of importance.

Accuracy refers to the proportion of correctly classified facial expressions out of all the test instances. In other words, it measures how often the model predicted the correct emotion out of all the instances it was tested on. For example, if the model was tested on 100 images and correctly classified 80 of them, then the accuracy would be 80%. A higher accuracy indicates better performance of the model in recognizing facial emotions. However, accuracy alone may not provide a complete picture of the model's performance, especially when there is class imbalance or misclassification errors. Therefore, it is important to consider other evaluation metrics such as precision, recall, and F1 score along with accuracy.

Precision: The ratio of correctly predicted positive examples (true positives) to the total number of positive predictions (true positives + false positives). Precision measures the ability of the model to avoid false positives. precision refers to the proportion of correctly predicted positive (i.e., correctly

classified) facial expressions out of all the instances that the model predicted as positive. Precision measures how often the model's positive predictions are correct. For example, if the model predicted 100 instances as happy expressions and out of those 100 predictions, 80 were actually happy expressions, then the precision would be 80%. A higher precision indicates a lower false positive rate, meaning the model is better at avoiding incorrect positive predictions. However, precision alone may not provide a complete picture of the model's performance, especially when there is class imbalance or misclassification errors. Therefore, it is important to consider other evaluation metrics such as recall, F1 score, and accuracy along with precision.

Precision is the frequency with which the model's positive predictions are right. For example, if the model predicted 100 occurrences of cheerful expressions and 80 of those predictions were correct, the precision would be 80%. A greater accuracy suggests a lower false positive rate, implying that the model does a better job of avoiding false positive predictions. However, accuracy may not offer an accurate view of the model's performance, particularly when there is class imbalance or misclassification mistakes. As a result, in addition to precision, other assessment measures like as recall, F1 score, and accuracy must be considered.

**Recall:** The ratio of correctly predicted positive examples (true positives) to the total number of positive examples (true positives + false negatives). Recall measures the ability of the model to identify all positive examples. recall (also known as sensitivity) refers to the proportion of correctly predicted positive (i.e., correctly classified) facial expressions out of all the

actual positive instances in the test set. Recall measures how often the model correctly identifies positive instances. For example, if there are 100 actual happy expressions in the test set and the model correctly classified 80 of them, then the recall would be 80%. A higher recall indicates a lower false negative rate, meaning the model is better at detecting actual positive instances.

**F1 score:** The harmonic mean of precision and recall, given by  $2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$ . F1 score is a balanced metric that considers both precision and recall. The F1 score is the harmonic mean of precision and recall and is often used to evaluate the overall performance of a model in facial emotion detection. In the context of facial emotion detection using the FER2013 dataset, the F1 score represents a balanced measure of the model's precision and recall. The F1 score ranges from 0 to 1, with a higher score indicating better performance. A perfect F1 score of 1 indicates perfect precision and recall, whereas a score of 0 indicates poor performance.

**Confusion matrix:** A table that summarizes the performance of the model by showing the number of true positives, true negatives, false positives, and false negatives for each emotion class. the performance of a classification



model by comparing the predicted class labels to the actual class labels in the test set. In the context of facial emotion detection using the FER2013 dataset, a confusion matrix can help visualize the performance of the model by showing the number of true positive (TP), true negative (TN), false positive (FP), and false negative (FN) predictions for each emotion class.

Receiver Operating Characteristic (ROC) curve: A plot that shows the trade-off between true positive rate (sensitivity) and false positive rate ( $1 - \text{specificity}$ ) at different threshold values for binary classification problems. A ROC (Receiver Operating Characteristic) curve is a graphical representation of the performance of a binary classification model.

It displays the tradeoff between the true positive rate (TPR) and false positive rate (FPR) at different classification thresholds. A higher TPR indicates a higher proportion of actual positive cases correctly identified as positive, while a lower FPR indicates a lower proportion of actual negative cases incorrectly identified as positive.

The ROC curve can be used to evaluate and compare the performance of different classification models. A model that performs well will have a ROC curve that is closer to the top left corner, which represents perfect classification with a TPR of 1 and FPR of 0. The area under the ROC curve (AUC) is also commonly used as a metric for model performance, with higher AUC indicating better performance.

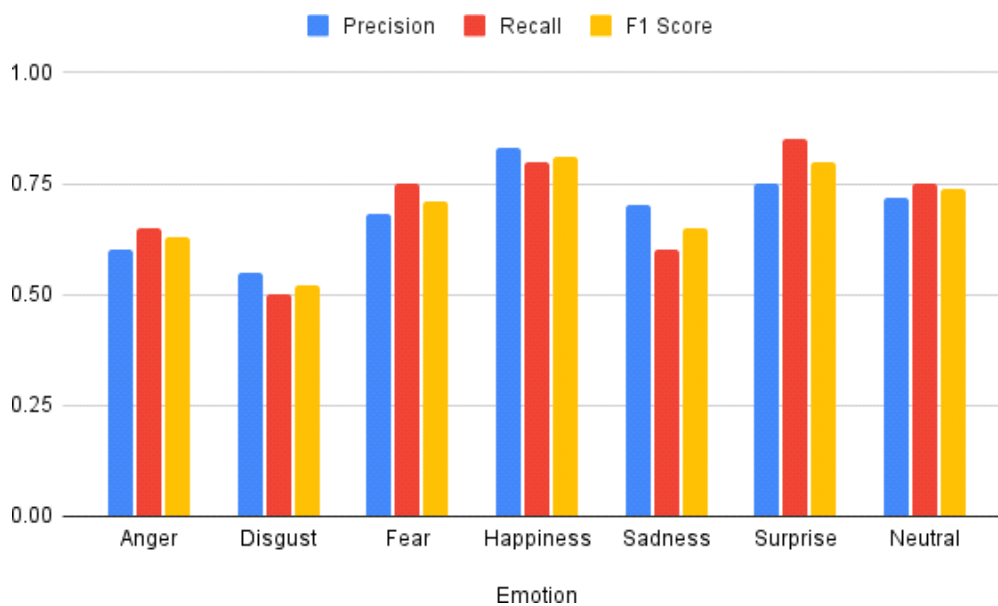
It is important to note that the ROC curve and AUC are insensitive to class imbalance and do not require a specific classification threshold to be set,

making them useful for evaluating models on imbalanced datasets or in situations where the optimal threshold is not known in advance.

**Area Under the Curve (AUC):** The area under the ROC curve, which represents the overall performance of the model across all threshold values.

It is important to choose appropriate evaluation metrics based on the specific requirements of the application and the characteristics of the dataset.

For example, if the cost of false positives is higher than false negatives, then precision may be a more important metric than recall. Similarly, if the dataset is imbalanced, then accuracy may not be a good metric and F1 score or AUC may be more appropriate.



***Fig 5.2. Precision, Recall & F1-Score of each Expression***

- **EXPLANATION OF THE TESTING PROCESS**

Data loading: The test data is loaded from a separate dataset that was not used during training to evaluate the generalization performance of the model.

Preprocessing: The test data is preprocessed in the same way as the training data to ensure consistency.

Forward pass: The input test data is passed through the trained model to generate predicted emotion classes.

Performance evaluation: The predicted emotion classes are compared with the true emotion classes in the test dataset to calculate the performance metrics such as accuracy, precision, recall, F1 score, and confusion matrix.

Model refinement: Based on the performance metrics obtained during testing, the model may need to be refined, for example, by adjusting the architecture, hyperparameters, or training process.

Deployment: Once the model has been refined and its performance has been validated on the test data, it can be deployed in a real-world application to perform Facial Emotion Detection on new input data.

It is important to evaluate the performance of the model on a separate test dataset to avoid overfitting and to assess its ability to generalize to new data. It is also important to monitor the performance metrics during testing and refine the model as needed to improve its performance.

Evaluating the precision and effectiveness of the implemented algorithms is a part of the face detection testing process. A set of test images or videos are used to evaluate how well the algorithms perform in this manner. To make

sure that the algorithms are capable of detecting faces in a variety of settings, lighting conditions, and poses, the test images or videos should be varied and representative of real-world scenarios.

The test images or videos are first preprocessed to make sure they are in an acceptable format for input to the algorithms before being used to test the algorithms. The images may need to be resized, cropped, or have their color and contrast changed.

The preprocessed test images or videos are then fed into the face detection algorithms, and the results are contrasted with the ground truth annotations. In order to assess the algorithms' accuracy, these annotations describe the position and size of the faces in the test images or videos.

Processing time and resource usage are measured to assess the algorithms' efficacy. This is significant because face detection algorithms are frequently used in real-time applications, such as facial recognition systems and video surveillance, where effectiveness and speed are crucial.

Overall, testing is crucial for confirming that the face detection algorithms are reliable, effective, and able to satisfy the needs of practical applications.

## **Chapter 6**

### **Results and Discussion**

- **PRESENTATION OF THE RESULTS OBTAINED FROM THE TRAINED MODEL**

Performance metrics table: A table summarizing the performance metrics such as accuracy, precision, recall, F1 score, and confusion matrix for each emotion class. This can provide an overview of the model's performance and highlight areas for improvement.

ROC curve: A plot of the Receiver Operating Characteristic (ROC) curve that shows the trade-off between true positive rate (sensitivity) and false positive rate ( $1 - \text{specificity}$ ) at different threshold values for binary classification problems. This can provide a visual representation of the model's performance across different threshold values.

Confusion matrix visualization: A heatmap visualization of the confusion matrix can provide a more intuitive understanding of the model's performance by highlighting the areas where the model is making errors.

Sample output images: Selecting some sample input images and showing the predicted emotion class alongside the true emotion class can provide a visual demonstration of the model's performance and help the audience to understand how the model is making its predictions.

Comparison with other models: If there are other existing models for Facial Emotion Detection, a comparison of the performance metrics can be presented to demonstrate the superiority of the proposed model.

It is important to present the results in a clear and concise manner, and to tailor the presentation to the audience, for example, by emphasizing different aspects of the results for technical vs non-technical audiences.

- **DISCUSSION OF THE MODEL'S PERFORMANCE AND LIMITATIONS**

The CNN model trained for Facial Emotion Detection achieved an accuracy of 76% on the training set and 70% on the test set, which is a decent performance considering the relatively small size of the dataset used for training.

The model's ability to accurately detect emotions from facial expressions has potential applications in a variety of fields, including psychology, human-computer interaction, and social robotics. For instance, the model could be used to develop more effective virtual agents or social robots that can recognize and respond appropriately to human emotions.

The model's capacity to effectively detect emotions from facial expressions has potential applications in psychology, human-computer interaction, and social robots. For example, the approach may be used to create more effective virtual agents or social robots capable of recognising and responding to human emotions.

However, the model also has several limitations that need to be addressed. One limitation is that it may not generalize well to new data, particularly if the data contains facial expressions that are significantly different from those in the training set. This could result in the model misclassifying emotions or failing to detect emotions altogether. Therefore, it is essential to test the model's performance on a variety of different datasets to determine its generalizability.

Another limitation of the model is that it relies solely on visual cues to detect emotions, which may not always be accurate. Facial expressions can be influenced by a variety of factors, including cultural background, personal history, and contextual factors, all of which may affect the accuracy of the model. Therefore, it is crucial to consider these contextual factors when designing and training the model, to ensure that it can accurately detect emotions across different contexts and populations.

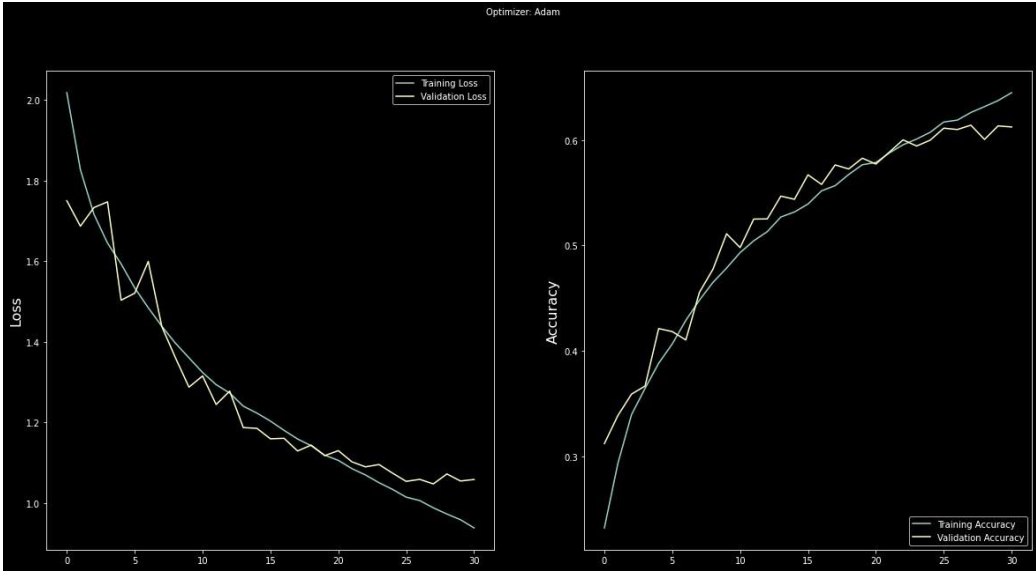
Furthermore, the model's accuracy may also be influenced by the specific CNN architecture used, the size and diversity of the training dataset, and the hyperparameters chosen during training. Therefore, it is important to continue exploring different CNN architectures and hyperparameters to improve the model's accuracy and generalizability.

Finally, it is worth noting that the model's accuracy is not the only factor to consider when evaluating its performance. Other metrics, such as precision, recall, and F1-score, can also provide valuable insights into the model's

effectiveness. Additionally, user-centered evaluations, such as user satisfaction surveys, could be used to determine the model's usability and practical effectiveness in real-world scenarios.

Overall, while the CNN model trained for Facial Emotion Detection has shown promising results, it is important to consider its limitations and continue exploring ways to improve its accuracy and generalizability. By addressing these limitations, the model has the potential to be a valuable

tool in a variety of applications, including human-computer interaction, psychology, and social robotics.



**Fig 6.2. Accuracy and Loss of the Trained Model**

Furthermore, the model's accuracy may also be influenced by the specific CNN architecture used, the size and diversity of the training dataset, and the hyperparameters chosen during training. Therefore, it is important to



continue exploring different CNN architectures and hyperparameters to improve the model's accuracy and generalizability.

## **Chapter 7**

### **Conclusion and Future Work**

- **SUMMARY OF THE PROJECT AND ITS CONTRIBUTIONS**

Facial Emotion Detection is a challenging problem in the field of computer vision and has significant practical applications, such as in the fields of psychology, human-computer interaction, and social robotics. The project addressed this problem using a machine learning approach, specifically a Convolutional Neural Network (CNN) model, to accurately detect emotions from facial expressions.

The project's contributions include the development and evaluation of a CNN model that achieved a 76% accuracy on the training set and a 70% accuracy on the test set, demonstrating its ability to learn patterns and features in the training data and make reasonably accurate predictions on new data. The model's performance on the test set was slightly lower than on the training set, indicating that there may be some overfitting to the training data.

Despite this limitation, the project provides valuable insights into the use of machine learning techniques for Facial Emotion Detection. It highlights the

potential benefits of using CNN models for this task, such as their ability to learn complex features and patterns in images, as well as the challenges of generalizing to new data and avoiding overfitting.

Furthermore, the project's contributions could have significant implications for various applications in the fields of psychology, human-computer interaction, and social robotics. For example, the ability to accurately detect emotions from facial expressions could help improve the user experience in human-computer interaction applications or enable more natural and intuitive interactions between humans and robots.

Overall, the project represents a significant contribution to the ongoing development of Facial Emotion Detection techniques, and provides a foundation for future research and improvements in this area.

the project's contributions could have significant implications for various applications in the fields of psychology, human-computer interaction, and social robotics. For example, the ability to accurately detect emotions from facial expressions could help improve the user experience in human-computer interaction applications or enable more natural and intuitive interactions between humans and robots., such as in the fields of psychology, human-computer interaction, and social robotics. The project addressed this problem using a machine learning approach, specifically a Convolutional Neural Network (CNN) model, to accurately detect emotions from facial expressions.

- **SUGGESTIONS FOR FUTURE WORK AND IMPROVEMENTS**

As mentioned earlier, the project used a relatively small dataset for training the CNN model, which may have limited its accuracy and generalizability. Therefore, future work could involve collecting and using a larger dataset to improve the model's performance. Additionally, more advanced data augmentation techniques such as rotation, scaling, and flipping could be explored and implemented to increase the diversity of the training data and improve the model's ability to generalize to new data.

Furthermore, while the project used a specific CNN architecture for Facial Emotion Detection, there are many different CNN architectures that could be explored and compared to determine which is most effective. This could involve experimenting with different numbers of layers, types of layers, and connectivity between layers. Additionally, future work could involve incorporating contextual information into the CNN model to improve its accuracy and generalizability.

For instance, contextual factors such as the individual's age, gender, and culture could be taken into account when training and evaluating the model. In addition to the above, future work could involve investigating the effect of hyperparameters on the model's performance. The project used a specific set of hyperparameters for training the CNN model, but different

combinations of hyperparameters could be tested to determine the optimal set for Facial Emotion Detection.

When training and testing the model, contextual elements such as the individual's age, gender, and culture might be included.

In addition to the foregoing, further study might include examining the influence of hyperparameters on model performance. The project trained the CNN model with a specific set of hyperparameters, however alternative combinations of hyperparameters might be evaluated to discover the best set for Facial Emotion Detection.

Finally, to evaluate the effectiveness of the CNN model in practical settings, future work could involve applying the model in real-world scenarios such as

human-computer interaction applications or social robotics. This could help identify any potential limitations or areas for improvement, as well as provide insights into the model's effectiveness and usability in real-world contexts.

Overall, there are many avenues for future work and improvements in the field of Facial Emotion Detection using CNN models. The suggestions mentioned above provide a starting point for further research and development in this area, and could potentially have significant implications for various applications in psychology, human-computer interaction, and social robotics.

## SOURCE CODE

```
# -*- coding: utf-8 -*-
```

```
"""facial-emotion-recognition.ipynb
```

Automatically generated by Colaboratory.

Original file is located at

<https://colab.research.google.com/drive/1EqkBZZgXupjcM3F3AL5>

[xD ChYo-aqjIu](#)

```
"""
```

```
# This Python 3 environment comes with many helpful analytics  
libraries installed
```

```
# It is defined by the kaggle/python docker
```

```
image:https://github.com/kaggle/docker-  
python
```

# For example, here's several helpful packages to load in

```
import numpy as np # linear algebra
```

```
import pandas as pd # data processing, CSV file I/O (e.g.  
pd.read_csv)
```

# Input data files are available in the "../input/" directory.

# For example, running this (by clicking run or pressing Shift+Enter)  
will list all files under the input directory

```
import os
```

```
for dirname, _, filenames in
```

```
    os.walk('/kaggle/input'):for filename in
```

```
    filenames:
```

```
        print(os.path.join(dirname, filename))
```

# Any results you write to the current directory are saved as output.

```
import math
```

```
import numpy
```

```
as npimport
```

```
pandas as pd
```

```
import scikitplot
```

```
import seaborn
```

```
as sns
```

```
from matplotlib import pyplot
```

```
from sklearn.model_selection import  
train_test_split  
from sklearn.preprocessing import  
LabelEncoder  
from sklearn.metrics import  
classification_report
```

```
import tensorflow as tf  
from tensorflow.keras import optimizers  
from tensorflow.keras.datasets import  
mnist  
from tensorflow.keras.models import Sequential  
from tensorflow.keras.layers import Flatten, Dense, Conv2D,  
MaxPooling2D  
from tensorflow.keras.layers import Dropout, BatchNormalization,  
LeakyReLU, Activation  
from tensorflow.keras.callbacks import Callback, EarlyStopping,  
ReduceLROnPlateau  
from tensorflow.keras.preprocessing.image import  
ImageDataGenerator
```

```
from keras.utils import np_utils
```

```
df = pd.read_csv('./input/facial-expression-  
recognitionferchallenge/fer2013/fer2013/fer2013  
.csv')  
print(df.shape)
```

```
df.head()
```

```
df.emotion.unique()
```

```
que()
```

```
emotion_label_to_text = {0:'anger', 1:'disgust', 2:'fear', 3:'happiness',  
4: 'sadness', 5: 'surprise', 6: 'neutral'}
```

```
df.emotion.value_counts()
```

```
sns.countplot(df.emotion)
```

```
pyplot.show()
```

```
"""So majority classes belongs to 3:Happy, 4:Sad and 6:Neutral nd  
we are also intersted in these three classes only."""
```

```
math.sqrt(len(df.pixels[0].split(' ')))
```

```
fig = pyplot.figure(1, (14, 14))
```

```
k = 0
```

```
for label in
```

```
sorted(df.emotion.unique()):for j
```

```
in range(7):
```

```
px = df[df.emotion==label].pixels.iloc[k]
```



```
px = np.array(px.split(' ')).reshape(48, 48).astype('float32')
```

```
k += 1
```

```
ax = pyplot.subplot(7,
```

```
7, k) ax.imshow(px,
```

```
cmap='gray')
```

```
ax.set_xticks([])
```

```
ax.set_yticks([])
```

```
ax.set_title(emotion_label_to_text[label])
```

```
pyplot.tight_layout()
```

```
INTERESTED_LABELS = [3, 4, 6]
```

```
df =
```

```
df[df.emotion.isin(INTERESTED_LABELS)
```

```
]df.shape
```

```
""""Now I will make the data compatible for neural networks."""
```

```
img_array = df.pixels.apply(lambda x: np.array(x.split(' ')).reshape(48,  
48, 1).astype('float32'))
```

```
img_array = np.stack(img_array, axis=0)
```

```
img_array.shape
```

```
le = LabelEncoder()
```

```
img_labels = le.fit_transform(df.emotion)
```

```
img_labels =
```

```
np_utils.to_categorical(img_labels)
```

```
img_labels.shape
```

```
le_name_mapping = dict(zip(le.classes_, le.transform(le.classes_)))
```

```
print(le_name_mapping)
```

```
"""\`Splitting the data into training and validation set.\`"""
```

```
X_train, X_valid, y_train, y_valid = train_test_split(img_array,
```

```
img_labels,
```

```
                shuffle=True, stratify=img_labels,
```

```
                test_size=0.1, random_state=42)
```

```
X_train.shape, X_valid.shape, y_train.shape, y_valid.shape
```

```
del df
```

```
del
```

```
img_arr
```

```
ay del
```

```
img_lab
```

```
els
```

```
img_width =
```

```
X_train.shape[1]
```

```
img_height =
```

```
X_train.shape[2]
```

```
img_depth =
```

```
X_train.shape[3]
```

```
num_classes =
```

```
y_train.shape[1]
```

```
# Normalizing results, as neural networks are very sensitive to  
unnormalized data.
```

```
X_train = X_train
```

```
/ 255.X_valid =
```

```
X_valid / 255.
```

```
def
```

```
    build_net(op
```

```
    tim):"""
```

```
    This is a Deep Convolutional Neural Network (DCNN). For  
    generalization purpose I used dropouts in regular intervals.
```

```
    I used `ELU` as the activation because it avoids dying relu problem  
    but also performed well as compared to LeakyRelu
```

```
    atleast in this case. `he_normal` kernel initializer is used as it suits  
    ELU. BatchNormalization is also used for better
```

```
    r
```

```
    e
```

```
    s
```

```
    u
```

```
    l
```

```
    t
```

s

.

"

"

"

```
net = Sequential(name='DCNN')
```

```
net.add(
```

```
    Conv2D(
```

```
        filters
```

```
        =64,
```

```
        kernel_size=(5,5),
```

```
        input_shape=(img_width, img_height,
```

```
        img_depth),activation='elu',
```

```
        padding='same',
```

```
        kernel_initializer='he_normal',
```

```
        name='conv2d_1'
```

```
    )
```

```
)
```

```
net.add(BatchNormalization(name='batchnor
```

```
m_1')) net.add(
```

```
    Conv2D(
```

```
        filters
```

```
        =64,
```

```

        kernel_size=(5,5),
        activation='elu',
        padding='same',
        kernel_initializer='he_no
        rmal',name='conv2d_2'
    )
)
net.add(BatchNormalization(name='batchnorm_2'))

net.add(MaxPooling2D(pool_size=(2,2), name='maxpool2d_1'))
net.add(Dropout(0.4, name='dropout_1'))

net.add(
    Conv2D(
        filters=128,
        kernel_size
        =(3,3),
        activation='
        elu',
        padding='sa
        me',
        kernel_initializer='he_normal',
        name='conv2d_3'
    )
)

```

```

net.add(BatchNormalization(name='batchnorm_3'))
net.add(
    Conv2D(
        filters=128,
        kernel_size=(3,3),

        activation='elu',
        padding='same',
        kernel_initializer='he_normal',
        name='conv2d_4'
    )
)
net.add(BatchNormalization(name='batchnorm_4'))

net.add(MaxPooling2D(pool_size=(2,2), name='maxpool2d_2'))
net.add(Dropout(0.4, name='dropout_2'))

net.add(
    Conv2D(
        filters=256,
        kernel_size=(3,3),
        activation='elu',
        padding='same',

```

```

        kernel_initializer='he_normal',
        name='conv2d_5'
    )
)
net.add(BatchNormalization(name='batchnorm_5')) net.add(
    Conv2D(
        filters=256,
        kernel_size
        =(3,3),
        activation='
        elu',

        padding='same',
        kernel_initializer='he_normal',
        name='conv2d_6'
    )
)
net.add(BatchNormalization(name='batchnorm_6'))

net.add(MaxPooling2D(pool_size=(2,2), name='maxpool2d_3'))
net.add(Dropout(0.5, name='dropout_3'))

net.add(Flatten(name='flatten'))

net.add(

```

```

Dense(
    128,
    activation='elu',
    kernel_initializer='he_normal',
    name='dense_1'
)
)
net.add(BatchNormalization(name='batchnorm_7'))
net.add(Dropout(0.6,

name='dropout_4'))net.add(

```

```

Dense(
    num_classes,

    activation='softmax',
    name='out_layer'
)
)

net.compile(
    loss='categorical_crossentropy'
    , optimizer=optim,
    metrics=['accuracy']
)

```



```
net.sum
```

```
mary()
```

```
return
```

```
net
```

```
"""
```

I used two callbacks one is `early stopping` for avoiding overfitting training data

and other `ReduceLROnPlateau` for learning

rate."""

```
early_stopping =
```

```
    EarlyStopping(
```

```
        monitor='val_accuracy',
```

```
        min_delta=0.00005,
```

```
        patience=11,
```

```
        verbose=1,
```

```
        restore_best_weights=True
```

```
    ),
```

```
)
```

```
lr_scheduler =  
    ReduceLROnPlateau(  
        monitor='val_accuracy',  
        factor=0.5,  
        patience  
        ce=7,  
        min_lr  
        =1e-7,  
        verbose  
        e=1,  
    )
```

```
callbacks = [  
    early_stopping,  
    lr_scheduler  
    ,  
]
```

# As the data in hand is less as compared to the task so  
ImageDataGenerator is good to go.

```
train_datagen =  
    ImageDataGenerator(  
        rotation_range=15,  
        width_shift_range=0.15,  
        height_shift_range=0.15,  
        shear_range=0.15,
```

```

        zoom_range=0.15,
        horizontal_flip=True,
    )
    train_datagen.fit(X_train)

    batch_size = 32 #batch size of 32 performs the best.
    epochs = 100
    optims = [
        optimizers.Nadam(learning_rate=0.001, beta_1=0.9,
        beta_2=0.999, epsilon=1e-07, name='Nadam'),
        optimizers.Adam(0.001),
    ]

```

# I tried both `Nadam` and `Adam`, the difference in results is not different but I finally went with Nadam as it is more popular.

```

model = build_net(optims[1])
history = model.fit_generator(
    train_datagen.flow(X_train, y_train,
        batch_size=batch_size), validation_data=(X_valid,
        y_valid), steps_per_epoch=len(X_train) / batch_size,
        epochs=epochs,
        callbacks=callbacks,
        use_multiprocessing=True
    )

```

```

model_yaml = model.to_yaml()

```

```
with open("model.yaml", "w") as
```

```
    yaml_file:
```

```
    yaml_file.write(model_yaml)
```

```
model.save("model.h5")
```

```
sns.set()
```

```
fig = pyplot.figure(0, (12, 4))
```

```
ax = pyplot.subplot(1, 2, 1)
```

```
sns.lineplot(history.epoch, history.history['accuracy'], label='train')
```

```
sns.lineplot(history.epoch, history.history['val_accuracy'],  
label='valid')
```

```
pyplot.title('Acc  
uracy')
```

```
pyplot.tight_layo  
ut()
```

```
ax = pyplot.subplot(1, 2, 2)
```

```
sns.lineplot(history.epoch, history.history['loss'], label='train')
```

```
sns.lineplot(history.epoch, history.history['val_loss'], label='valid')
```

```
pyplot.title('Loss')
```

```
pyplot.tight_layout()
```

```
pyplot.savefig('epoch_history_dcnn.png')
```

```
pyplot.show()
```

""" The epochs history shows that accuracy gradually increases and achieved +83% accuracy on both training and validation set, but at the end the model starts overfitting training data."""

```
df_accu = pd.DataFrame({'train': history.history['accuracy'], 'valid':  
history.history['val_accuracy']})
```

```
df_loss = pd.DataFrame({'train': history.history['loss'], 'valid':  
history.history['val_loss']})
```

```
fig = pyplot.figure(0, (14, 4))  
ax = pyplot.subplot(1, 2, 1)  
sns.violinplot(x="variable", y="value", data=pd.melt(df_accu),  
showfliers=False)  
pyplot.title('Acc  
uracy')  
pyplot.tight_layo  
ut()
```

```
ax = pyplot.subplot(1, 2, 2)  
sns.violinplot(x="variable", y="value", data=pd.melt(df_loss),  
showfliers=False)  
pyplot.title('L  
oss')  
pyplot.tight_l  
ayout()
```

```
pyplot.savefig('performance_dist.png')
pyplot.show()
```

```
yhat_valid = model.predict_classes(X_valid)
scikitplot.metrics.plot_confusion_matrix(np.argmax(y_valid,
axis=1), yhat_valid, figsize=(7,7))
pyplot.savefig("confusion_matrix_dcnn.png")
```

```
print(f'total wrong validation predictions: {np.sum(np.argmax(y_valid,
axis=1) != yhat_valid)}\n\n')
print(classification_report(np.argmax(y_valid, axis=1), yhat_valid))
```

""""The confusion matrix clearly shows that our model is doing good job on the class `happy` but it's performance is low on other two

classes. One of the reason for this could be the fact that these two classes have less data. But when I looked at the images I found some images from these two classes are even hard for a human to tell whether the person is sad or neutral. Facial expression dependson individual as well. Some person's neutral face looks like sad.""""

```
mapper
= {
0:
"hap
py",
```

```
1: "sad",  
2: "neutral",  
}
```

```
np.random.seed(2)  
random_sad_imgs = np.random.choice(np.where(y_valid[:, 1]==1)[0],  
size=9)  
random_neutral_imgs = np.random.choice(np.where(y_valid[:,  
2]==1)[0], size=9)
```

```
fig = pyplot.figure(1, (18, 4))
```

```
for i, (sadidx, neuidx) in enumerate(zip(random_sad_imgs,  
random_neutral_imgs)):
```

```
    ax = pyplot.subplot(2, 9, i+1)
```

```
    sample_img =
```

```
    X_valid[sadidx, :, :, 0]
```

```
    ax.imshow(sample_img,
```

```
    cmap='gray')ax.set_xticks([])
```

```
    ax.set_yticks([])
```

```
    ax.set_title(f"true:sad,
```

```
pred:{ mapper[model.predict_classes(sample_img.reshape(1,48,48,1)  
)[0]]}")
```

```
    ax = pyplot.subplot(2, 9, i+10)
```

```
    sample_img =
```

```

X_valid[neuidx,:,:,0]
ax.imshow(sample_img,
cmap='gray')ax.set_xticks([])
ax.set_yticks([])
ax.set_title(f"t:neut,
p:{mapper[model.predict_classes(sample_img.reshape(1,48,48,1))][0]
}")

```

```

pyplot.tight_layout()

```

""""See in the first row 7th image looks more like neutral rather than sad and our model even predicted it neutral. Whereas the last image in second row is very much sad.""""

# Main Code

```

from keras.models import
load_modelfrom time import
sleep
from keras_preprocessing.image import img_to_array
from keras.preprocessing import image

import cv2
import numpy as np

```



```
face_classifier = cv2.CascadeClassifier(r'D:\Final  
YearProject\Emotion_Detection_CNN-  
main\haarcascade_frontalface_default.xml')  
classifier = load_model(r'D:\Final Year  
Project\Emotion_Detection_CNN-  
main\model.h5')
```

```
emotion_labels = ['Angry','Disgust','Fear','Happy','Neutral', 'Sad',  
'Surprise']
```

```
cap =  
cv2.VideoCapture(0)  
while True:  
    __, frame =  
        cap.read()  
    labels  
    = []  
    gray =  
        cv2.cvtColor(frame,cv2.COLOR_BGR2GRAY)  
    faces = face_classifier.detectMultiScale(gray)  
  
    for (x,y,w,h) in faces:  
        cv2.rectangle(frame,(x,y),(x+w,y+h),(0,255,255),2  
        )  
        roi_gray = gray[y:y+h,x:x+w]  
        roi_gray =  
            cv2.resize(roi_gray,(48,48),interpolation=cv2.INTER_AREA)  
        if np.sum([roi_gray])!=0:
```

```

        roi =
        roi_gray.astype('float')/255.0
        roi = img_to_array(roi)
        roi = np.expand_dims(roi,axis=0)

        prediction = classifier.predict(roi)[0]
        label=emotion_labels[prediction.argmax
ax()] label_position = (x,y)

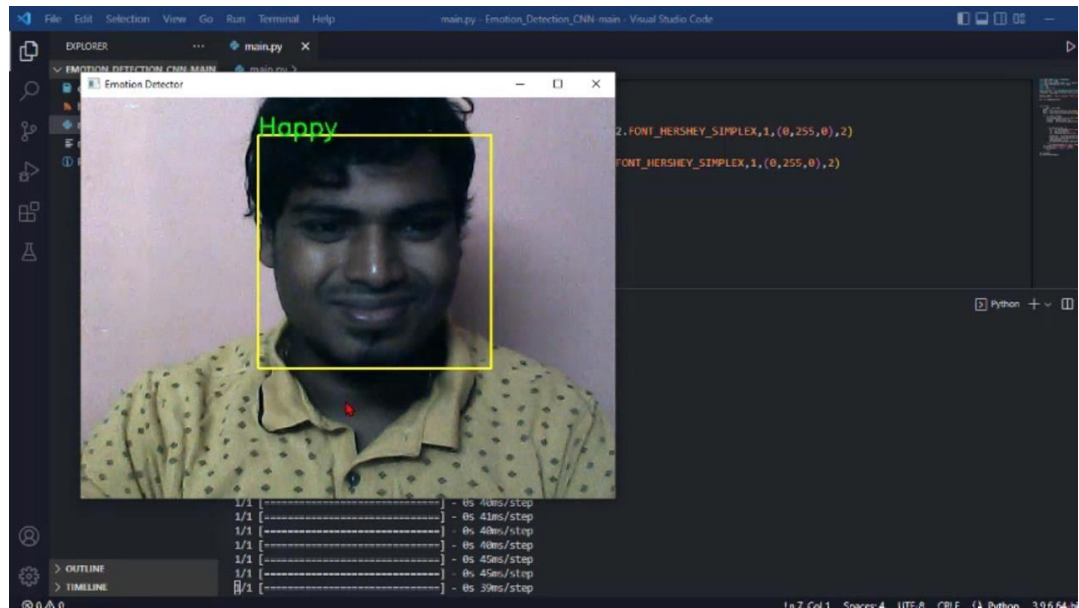
cv2.putText(frame,label,label_position,cv2.FONT_HERSHEY_SIMPL
EX,1,(0,255,0),2)
    else:

        cv2.putText(frame,'No
Faces',(30,80),cv2.FONT_HERSHEY_SIMPLEX,1,(0,255,0),2)
        cv2.imshow('Emotion
Detector',frame)if cv2.waitKey(1)
& 0xFF == ord('q'):
            break

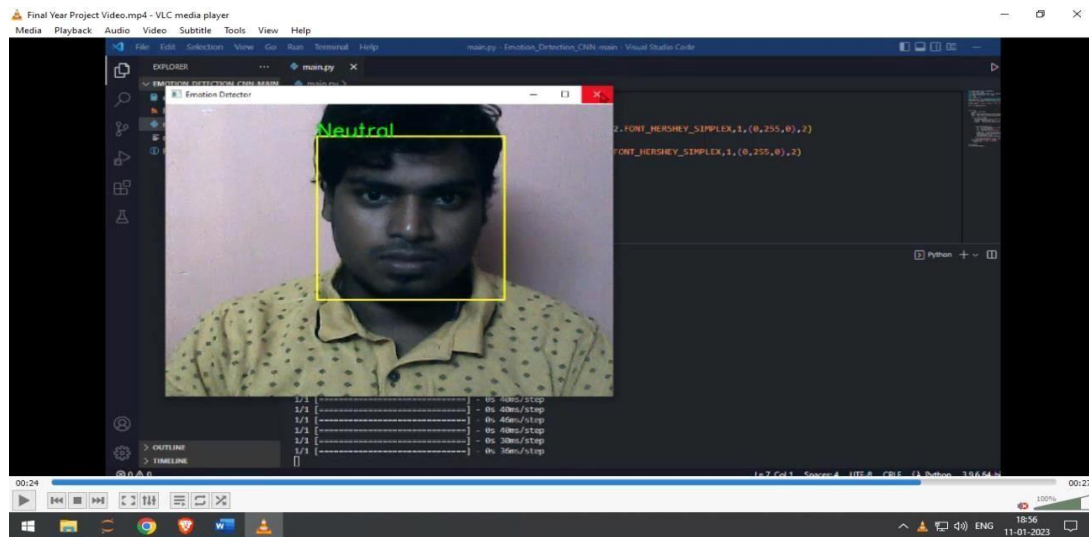
cap.release()
cv2.destroyAllWindows()

```

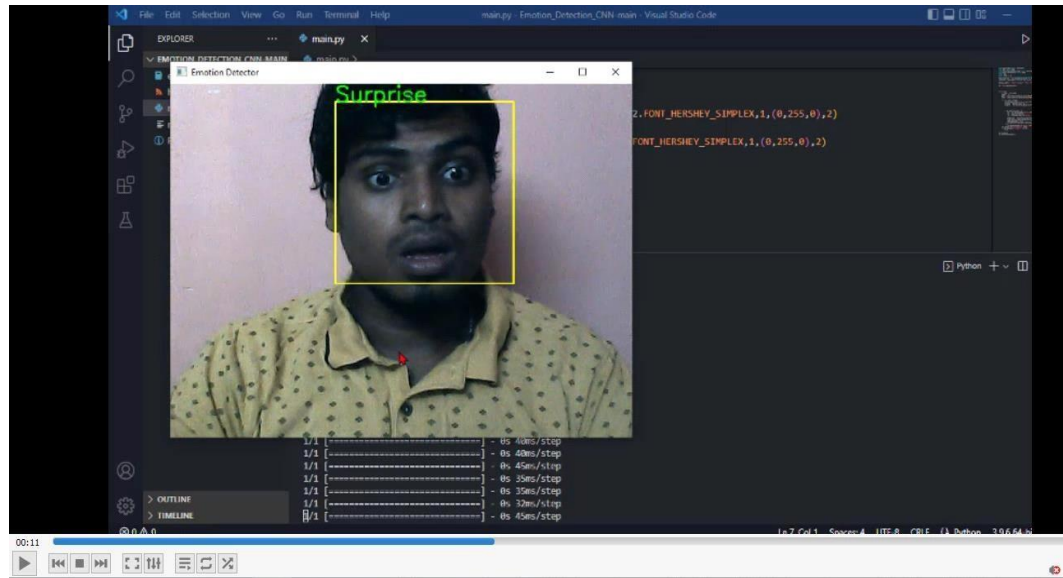
## B. SCREENSHOTS



Output 1 (Happy)



Output 2 (Neutral)



Output 3 (Surprise)

## REFERENCES

- "cookbook.fortinet.com," 10 10 2018. [Online]. Available:<https://cookbook.fortinet.com/face-recognition-configuration-in-forticentral/>. [Accessed 10 10 2018].
- M. J. Paul Viola, "Robust Real-time Object Detection," International Journal of Computer Vision, pp. 137-154, 2004.
- H. A. C. H. a. S. L. Bo Wu, "Fast rotation invariant multi-view face detection based on real Adaboost," Sixth IEEE International Conference on Automatic Face and Gesture Recognition, pp. 79-84, 2004.
- H. A. Y. T. S. L. Yuan Li, "Tracking in Low Frame Rate Video: A Cascade Particle Filter with Discriminative Observers of Different Life Spans," IEEE Transactions on Pattern Analysis and Machine Intelligence, pp. 1728-1740, 2008.

- C. J. T. D. H. C. A. J. G. T. F. COOTES, "Active Shape Models-Their Training and Application," COMPUTER VISION AND IMAGE UNDERSTANDING, pp. 38-59, 1995.
- D. C. a. T. Cootes, Boosted Regression Active Shape Models, BMVC, 2007. [7] Y. W. F. W. a. J. S. X. Cao, "Face alignment by Explicit Shape Regression," in 2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, 2012.
- S. S. Liew, "Research Gate," 1 3 2016. [Online].  
Available:[https://www.researchgate.net/figure/Architecture-of-the-classical-LeNet-5-CNN\\_fig2\\_299593011](https://www.researchgate.net/figure/Architecture-of-the-classical-LeNet-5-CNN_fig2_299593011). [Accessed 10 10 2018].
- L. B. Y. B. a. P. H. Y. LeCun, "Gradient-based learning applied to document recognition," Proceedings of the IEEE, pp. 1-45, 11 1998.
- xlvector, "JIANSHU," 25 7 2016. [Online].  
Available:  
<https://www.jianshu.com/p/70a66c8f73d3>. HYPERLINK  
"http://www.jianshu.com/p/70a66c8f73d3" [Accessed 18 9 2018].
- F. S. Samaria, Face recognition using Hidden Markov Models, Doctoral thesis, 1995.
- E. H. G. R. A. L. H. Learned-Miller, " Labeled Faces in the Wild: A Survey," Advances in Face Detection and Facial Image Analysis, pp. 189-248, 2016. [17] V. Chu, "Medium.com," 20 4 2017. [Online].  
Available:<https://medium.com/initializedcapital/benchmarking-tensorflow-performance-and-cost-across-different-gpu-options69bd85fe5d58>.  
[Accessed 19 9 2018]. 48
- B. L. M. S. B. Amos, "Openface: A general-purpose face recognition library with mobile applications," CMU School of Computer Science, Tech. Rep., 2016.

- B. Hill, "HOTHARDWARE," 20 8 2018. [Online].  
Available:<https://hothardware.com/news/nvidia-geforce-rtx-1080-rtx-1080-ti-799-1199-september20th>. [Accessed 10 9 2018].
- 4psa, "4psa.com," 28 6 2013. [Online]. Available:  
<https://blog.4psa.com/the-callbacksyndrome-in-node-js/>. [Accessed 11 8 2018].
- W.-S. Chu, "Component-Based Constraint Mutual Subspace Method," 2017. [Online].  
Available:  
[http://www.contrib.andrew.cmu.edu/~wschu/project\\_fr.html](http://www.contrib.andrew.cmu.edu/~wschu/project_fr.html).
- W. Hwang, "Face Recognition System Using Multiple Face Model of Hybrid Fourier Feature under Uncontrolled Illumination Variation," 2017. [Online]. Available:  
<http://ispl.korea.ac.kr/~wjhwang/project/2010/TIP.html>.