# Review on Human Face Detection Using Deep Learning

A. Srinidhi[#1], S. Bhavanisankari[*2]

[1]*Post Graduate, Electronics and Communication Engineering Department, Jerusalem College of Engineering, Chennai.*
[2]*Associate Professor, Electronics and Communication Engineering Department, Jerusalem College of Engineering, Chennai.*
[1]email id: srinidhianand4@gmail.com [orcid: 0009-0000-1988-5046], [2]email id:
bhavanisankari@jerusalemengg.ac.in

*Abstract*— **In the real world, the detection of human face objects is challenging and seems to be a complex task. Emotions play a vital role for every individual to understand their own self and others and honour it. Emotional intelligence (EI) is one's ability to perceive others to manage better relationships. Facial Detection is the primary key in EI and to achieve this potential, artificial intelligence helps in EI's needs. In today's world, domains like mentoring, photography, security systems, social media platforms, augmented reality and so on, use face detection as a tool for assessment and surveillance. The significance of this review lies in its multidisciplinary focus, bridging the gap between computer science and psychiatry. As soon as one must identify facial expressions, the process gets complex. The demand to provide a promising model to recognize faces with respect to computational cost, processing speed and execution time even in low-power devices has become more prominent in real time applications. By focusing on emotion detection and social behaviour analysis, these systems can offer clinicians novel ways to assess and monitor psychiatric patients, thus contributing to personalized treatment strategies. In this work, various deep learning algorithms are tested to detect the faces of different resolutions. MTCNN, also known as Multi-Task Cascaded CNN, specializes in detecting and learning features from facial images that relate to facial expressions. This proposed system aims to accurately detect faces in real time by extracting the feature points and feeding them to the model to recognize the landmarks as an efficient way for effective and precise facial expression identification, with potential uses in a number of fields, such as robotics, human-computer interface, sales, and healthcare. Hence, through experimental analysis and performance visualization, this work identifies the best deep learning model achieving better accuracy with feasible performance, speed and execution time on large datasets when compared to other models.**

*Keywords*–— **Facial Detection, Artificial Intelligence, MTCNN, Deep Learning, psychiatric applications, non-invasive monitoring.**

## 1. INTRODUCTION

Computer vision is a component of AI which processes digital sources like images, videos etc., for information. Among them, Facial Detection is an interesting area of Artificial intelligence, where research has been going for many years to earn stability in the recognition system. It is generally termed as the identification or verification of a person's face from the database using the detection of presence of their face. To train the model with the features, the facial features are measured and retrieved from the face. Computers can detect and even comprehend human emotions thanks to facial recognition, a biometric approach that expresses and analyses human face patterns.

One of the modern uses of artificial intelligence (AI) with neural networks is the recognition of faces in images and videos for a range of purposes. The majority of methods process visual data and look for broad patterns in the faces of people in pictures or videos. Also, face detection usage is becoming increasingly popular for safety and security purposes in low-powered systems such as mobiles. Generally, AI algorithms techniques play a vital role in today's world. The algorithmic approach can give instantaneous results with the smart approach of searching the databases of faces to identify the emotion detected in an input.

The backbone of AI is Artificial Neural Networks which mimic the human biological neural system to track the recognition and detection activities. The discovery and implementation of methods for interpreting, encoding, and extracting characteristics from facial expressions is the aim of this field of study in order to improve computer prediction. Because of deep learning's remarkable success, performance is being improved by making the most of all of its different architectures. This work is motivated to analyse the precise CNN model network to detect faces objects in real-time at compromising time.

## 2. RESEARCH ON FACE DETECTION SYSTEM

A facial recognition system has been in use since the 1960s. It uses technology to compare a digital image or video frame with a database of known faces to identify a person's face. By recognizing and

measuring face features from the provided pictures using ID verification services, this system is used for user authentication.

### A. Deep Learning

Deep Learning is a branch of machine learning that uses algorithms inspired by neural networks or the way the human brain functions. We refer to these configurations as neural networks. The computer is trained to perform tasks that come naturally to people. Artificial Neural Networks (ANN), Recurrent Neural Networks (RNN), and reinforcement neural networks are some of the models that are utilized in deep learning. However, the Convolutional Neural Networks (CNN) has made necessary contributions to the computer vision and image analysis field.

Artificial neurons in different layers make up convolutional neural networks (CNNs). Artificial neurons are computational functions that provide an activation value as an output after adding up a variety of inputs. They perform comparable functions to the neuron cells that the brain utilizes to transmit different input messages from senses and other reactions. The values of each CNN neuron's weights dictate how it acts. When given pixel values, CNN's artificial neurons can recognize many visual features and characteristics.

CNN is helpful for image identification since it provides high accuracy. Image recognition is the first step in detecting faces and has many applications in the phone, security, medical picture analysis, and recommendation systems, among other areas. Many inbuilt libraries are provided by the Python language to implement this architecture. In brief, facial detection is proposed in four stages namely,

- Residual Blocks - Grid Cells or Anchor boxes or Kernel of any size is employed to divide the image into small boxes.
- Bounding box regression - Regressor module is used to derive the box coordinates to localize the faces in the input images.
- Intersection of Union - An estimation metric used to gauge an object detector's precision on a given dataset.
- Non Max Suppression - Selecting the bounding box of high score, above the provided threshold and eliminating other bounding boxes calculated over the images.
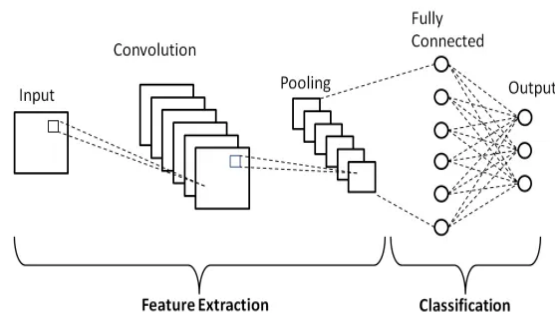


Fig. 1  A CNN model outline.

The convolutional neural network consists of two layers such as feature extraction layer and classification layer as shown via Fig 1. Feature Extraction part includes the input layer, Convolution layers and Pooling layers. The Classification part consists of a FC layer to detect the object class.

### B. Short Analysis on Previous Approaches

Face recognition has been the subject of research for a very long time. In order to present the Fast-FaceNet technique, some researchers have drawn comparisons between the FaceNet concept and

MobileNet architecture [1]. To improve the local feature extraction of face emotions, the attention module was first added to the MobileNetV1 model [3]. The model parameters are then adjusted by combining the centre loss and softmax loss in order to reduce intra-class distance and boost inter-class distance. Without adding more model parameters, the suggested approach greatly increases recognition accuracy when compared to the original MobileNet series models.

Furthermore, in order to extract the performance of facial feature algorithms, MTCNN learned a mapping from face images to a compact Euclidean space, whose distances relate directly to a measure of face similarity. This mapping was learned using Google's FaceNet framework. Additional research should be done on enhancing automatic cluster selection, particularly on increasing clustering accuracy through k-means, MTCNN, and hyper parameter modification [16].

The research uses a combination of MobileNet V2 and SSD to achieve excellent recognition accuracy and real-time expression recognition. The deep neural network can automatically extract the visual feature for precise classification because the robot's processing performance is constrained [2]. From the previous research works, MTCNN proved to give better results for the face detection and landmark position alignment when compared with other methods. It focuses on face detection in contrast to other object detection algorithms and achieves better accuracy, besides the lack of recognition tasks and its network computation complexities.

*C.  Face Detection*

The term "facial detection" refers to the capacity to locate a face in any input image or frame. The bounding box coordinates of the faces that have been detected, along with a confidence level ensuring that the face is localized in any input image or frame, are the output. The best component for processing images is Convolutional Neural Network (CNN), which also functions as the best algorithm for deep learning techniques.

The target detection algorithm has evolved from the conventional algorithm based on manual features to the detection technology based on deep neural networks in recent years due to the rapid growth of deep learning technology. One-stage and two-stage target identification algorithms are currently available as deep learning based methods.

*D.  Stages of Face Detection*

Generally, Object detection consists of either single or two stages. Detection architecture comprises of two types namely,

- Single stage detector.
- Two stage detector.

Appearance based face detection is implemented using neural networks. In traditional networks, the algorithms are based on the sliding window concept, have high computational complexity and limited detection effect.

As in Fig 2, Instead of doing regional proposals, one-stage target detection evenly performs dense sampling at many image regions. The regional convolution neural network series serves as the primary basis for the two-stage detection model.
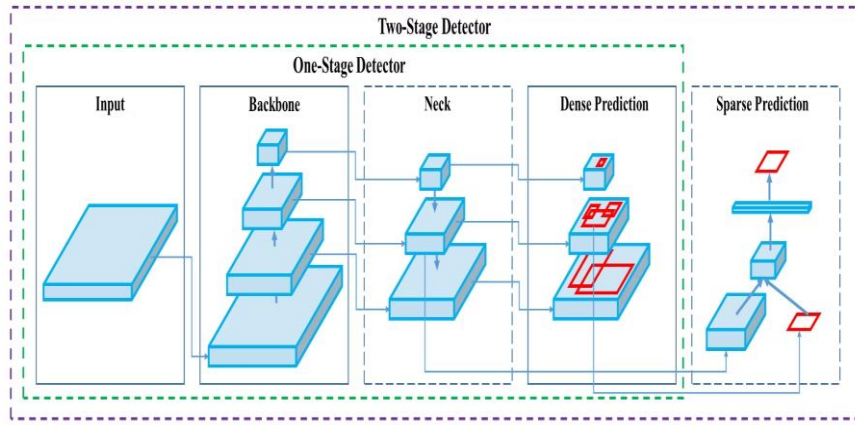
Fig. 2.  Face Detection Algorithm Stages.

MTCNN is based on the multi stages Convolution neural networks for face detection and face landmark detection like eyes, nose and mouth. Let's discuss the datasets used and methods experimented for face detection tasks.

This work explores various deep neural architectures and addresses the performance issues in analysing and detections of various emotions in real-time. They provide solutions on the development of a reduced architecture model and the improvement of computational accuracy.

In the following, section III analyses the various face detection methodologies implementation for detecting the inputs from the acquired images and performance metrics. Section IV Shows the experiment results and discussion related to training and testing datasets and its inferences. The section V Narrates the conclusion and future works.

## 3. METHODOLOGY IMPLEMENTATION

This paper explores different face detection algorithms on the defined datasets. Fig. 3 provides an overview of the proposed face detection technique.
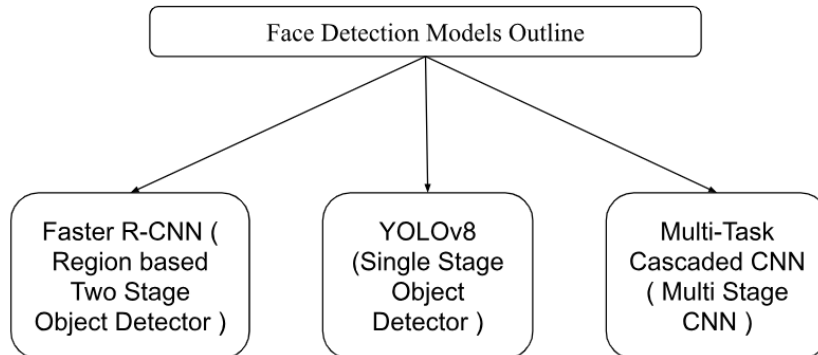


Fig. 3.  Face detection Model Outline.

The above described object detection models are built on the backbone neural architecture of various types. The neural network outline information is indicated in Table I. The object detection algorithm built on CSP DarkNet-53 and Inception ResNet v2 neural architectures. The Data augmentation techniques are also used in the YOLOv8 model.

TABLE I
NEURAL NETWORK OUTLINE STRUCTURE OF DETECTION ALGORITHM

| Faster R-CNN | YOLOv8 | MTCNN |
|---|---|---|
| **Backbone**: Inception ResNet v2 network with 164 layers<br><br>**Output**: Detects Bounding box at two stages.<br><br>Concentrates on 1000 classes by default | **Backbone**: CSP DarkNet 53 consist of 53 layers<br><br>**Neck:** Feature Pyramid Network + Path Aggregation Network<br><br>**Output**: Detects & classifies bounding box at single stage<br><br>Concentrates on 80 classes by default. | CNN with three layers<br>Proposal Network<br>Region Network<br>Output Network along with Face Landmark (eyes, nose, mouth)<br><br>**Output**:<br>Face Detection<br>Bounding Box<br>Landmark Detection |

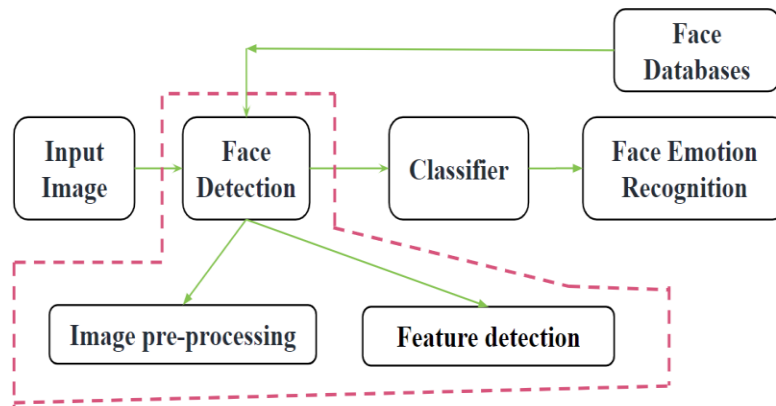The flow diagram for the entire process of FER is shown in Fig 4.



Fig. 4.  The System Block of Face Recognition.

The process starts with preparing the custom dataset of digital images of different resolutions from the various platforms. Then the images are processed using the OpenCV library functions to prepare the data shape to fit into the model. Model is trained on the datasets to output feature maps that are formed to detect and classify the object with the desired results. The proposed work comprises the face detection part, the foremost essential part in the FER system

*A. Faster R-CNN Method*

The algorithms for R-CNN and Fast R-CNN are mostly based on selective search; however, Faster R-CNN considerably accelerates the algorithm by substituting the region proposal network (RPN) for the selective search technique. Unlike the filter concatenation step of the Inception architecture, Inception-ResNet-v2 makes use of residual connections in a convolutional neural network. The Inception architecture is combined with residual connections to create Inception-ResNet. Fig. 5 illustrates the steps in the Faster R-CNN algorithm's working flow for object detection in an image.

Step 1: Provide images as input to ConvNet, and it will return feature maps for the image.
Step 2: Use these feature maps to apply Region Proposal Network (RPN) and obtain object proposals.
Step 3: Reduce each proposal's size to the same level by using the ROI pooling layer.

Step 4: In order to classify any predictions for the bounding boxes of the image, transmit these recommendations to a fully connected layer at the end.
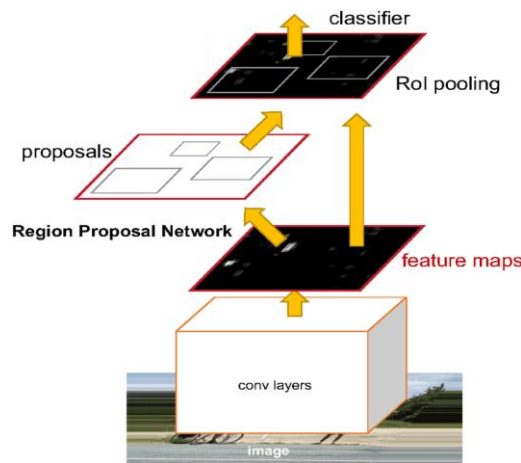


Fig. 5. The System Block of Face Recognition.

The disadvantage of this approach is that it takes time to suggest an object, and as systems run one after the other, the effectiveness of one system depends on the success of the system that comes before it.

*B. YOLOv8 Method*

The most recent model in the YOLO model series is called YOLOv8. YOLO models have the advantage of being faster to train and capable of producing acceptable accuracy at lower model sizes. In order for YOLOv8 to function, the input image is divided into a grid, usually measuring $13 \times 13$ or $26 \times 26$, depending on the kind. Each grid cell is responsible for anticipating items within its own spatial area.

Compared to its predecessors, YOLOv8 has more features, which include
1. In order to address the network's gradient expansion or disappearance issues, the residual layer is implemented.
2. The network for feature extraction that uses the input images to generate feature maps.
3. In addition to saving computing costs and network memory consumption, CSPNet improves the accuracy and speed of inference.

The basic steps of YOLOv8's operating principle are as follows:

1. YOLOv8 separates an image into a grid, usually measuring 13x13 or 26x26, depending on the kind of input called as input processing

2. A deep convolutional neural network (CNN) is used by the network to extract high-level features from the input image. Popular models like ResNet or Darknet are often used as a basis for network architecture decisions to this feature extraction level

3. By regressing the box's width, height, and top-left corner coordinates, YOLOv8 predicts an object's bounding box. Additionally, a confidence score is calculated to indicate the probability that the object is in the expected box and it's simply denoted as bounding box prediction

4. YOLOv8 predicts class probabilities for every grid cell in addition to bounding box predictions. This suggests that the model is capable of more than simply object detection; it can also identify objects and the categories that go along with them.

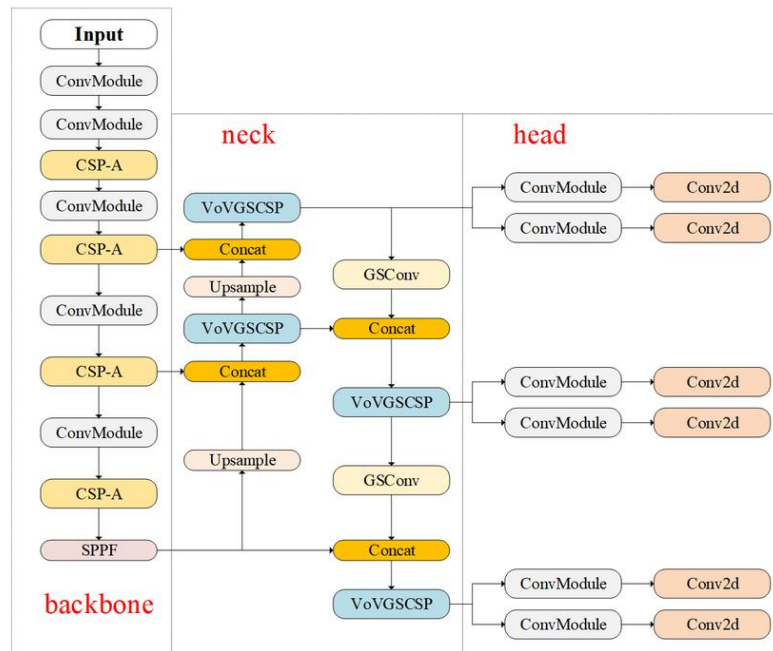Fig 6 illustrates the flow diagram of the YOLOv8 model approach.

Fig.6. YOLOv8 architecture

After the predictions are generated, low-confidence detections are filtered out by applying a confidence threshold. Non-maximum suppression is used to remove the overlapping or duplicate bounding boxes, ensuring that the most accurate search is the only one left running.

*C. Multi-task Cascaded CNN Method*

A cascade series of convolutional neural networks (CNNs) is used by the MTCNN (Multi-Task Cascaded Convolutional Neural Networks) algorithm, a deep learning based face recognition and matching method, in order to recognize and label faces in digital photos or videos. Three components make up Multi-task cascaded convolution networks (MTCNN), which are primarily employed for forecasting and detecting tasks. It falls under the system's operation. The fig. 7 shows the working flow which includes

- Stage 1: Proposal Network:

The bounding box regression vectors for the candidate windows are obtained using this proposal network stage. After obtaining the bounding box vector, the overlapping regions are merged by refining. Following the modification, fewer candidates will be considered.

- Stage 2: Refine Network

This Refine Network employs non-maximum suppression (NMS) to add overlapping candidates, calibrates its bounding box regression, and constantly decreases the number of candidates. R-Net produces a 4-element vector that serves as the face's bounding box and the face itself, regardless of whether the input is a face object or not. Ten element vectors for localizing landmarks.

- Stage 3: Output Network:

This Output network part is comparable to the second network stage; however, this output network's goal is to provide a detailed description of the face and to establish the location of five facial markers for the mouth, nose, and both left and right eyes endpoints.

Three different objectives are assigned to the network: one is facial landmark localization (e.g., mouth, nose, and eyes), followed by bounding box regression, and final one is face/non-face categorization.
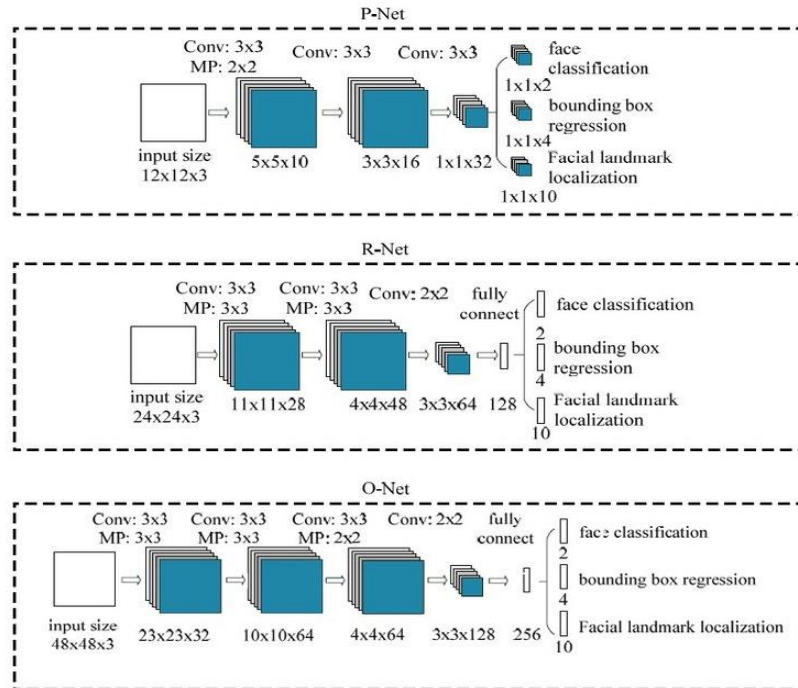
Fig.7  MTCNN Architecture

*D. Model Development using Python*

The environment setup for the pre-trained model made using programming language, Python. The following libraries used to develop the model.

*1)* *Numpy*: This library is for scientific computing in Python. They provide various functions for numerical operations like mean, average and filtering.

*2)* *Matplotlib:* This library is utilized to create interactive, animated, and static Python visualizations. Matplotlib allows for the visualization of both difficult and easy data.

*3)* *Ultralytics (for YOLO):* The deep learning package for loading the YOLO model for various tasks like predict, detect, classify the images or video or any input stream files to detect & classify the objects.

*4)* *MTCNN:* The deep learning CNN package for loading the MTCNN network to detect faces using the built in methods to refine the images and output the features with classification.

*5)* *Tensorflow:* One of python's machine learning libraries, utilized for loading the pretrained model for face detection and classification tasks.

*6)* *OpenCV:* One of the libraries for image manipulations like read, show, and write. Also, it supports real time acquisition of images using webcam.

*7)* *LabelImg:* Python based graphical representation of image annotation tool for training the YOLO models.  Annotations stored as text files in YOLO format. Besides, it also supports PASCAL VOC and CreateML formats.

The following steps are used to load the pretrained model and train and model is evaluated based on the custom dataset.
- Model is imported.
- Model is loaded using Tensorflow with default weights.

- Prepare the dataset for train & test.
- Compile the model.
- Fit the model with custom datasets.
- Validate the model with test Datasets.

Visualize the model's Confidence (or fitness or objectness score) along with the inference time measures using matplotlib package.

### E. Performance Metrics

The detector confidence (or objectness score) helps to evaluate its algorithm performance as one of the metrics. In general, confidence is calculated while detecting and generating bounding boxes. Confidence, represented as a percentage, is the likelihood that the image will be accurately identified by the algorithm.

The average total accuracy scores across all thresholds are used to calculate confidence scores. The following provides the average precision (confidence) measurement:

$$\text{Confidence} = \frac{1}{|Thresholds|} = \sum_t \frac{TP}{TP + FP}$$

TP - True Positive - Face object present and detected.

FP - False Positive - Face object present and not detected.

## 4. EXPERIMENTAL RESULTS AND DISCUSSION

The proposed research employs Google Colab as the experimental debugging platform and programming language as Python, Tensorflow, MTCNN, and Ultralytics YOLO to train network model parameters, and some OpenCV image library functions to help create and show the face detection window.

### A. Dataset Collection and Processing

The custom dataset has been developed using Google images along with the real time image acquisition technique. A total of 125 photos are gathered and divided into training and testing datasets at each resolution in a 1:4 ratios. To train the model with the target object such as the human face, a training set consisting of more than 100 images per resolution is evaluated. The training images are resized based on the user defined resolution.

The test datasets are categorized in the range of five different resolutions such as 40 x 40, 60 x 60, 160 x 160, 224 x 224, and 512 x 512. The test datasets consist of 25 real time images per resolution.

Custom dataset is prepared using Google random images of human faces of various ages and real time people facial images under various resolutions of low to high range are shown in Table II.

TABLE II
DATASETS COLLECTION

| | |
|---|---|
| **Training Datasets** | 125 Face Images are obtained from various platforms |
| **Training Resolutions** | 4 Resolutions (40*40, 60*60, 160*160, 224*224) |
| **Test Datasets**: | 25 Real Time faces images per resolutions |
| **Testing Resolutions** | 5 Resolutions (40*40, 60*60, 160*160, 224* 224, 512*512) |

### B. Training Set and Test Set Inferences

Calculation & Observations made on the following aspects and inferences are observed.
- Overall Training Time (Inference Time) in seconds.

● Overall Objectness Confidence (Objectness Score) in percentage.
● Overall Test Image Predictions are analyzed.

The various tested object detection algorithm analyzed with pretrained weights & custom class namely
● YOLOv8
● Faster R-CNN
● MTCNN

The Features of each methodology is presented in Table III.

TABLE III
DATASET AND OUTPUT INFERENCE OF FACE DETECTION METHODS

| YOLOv8 | Fully concentrates on the entire bounded object to extract all features |
|---|---|
| Faster R-CNN | All Features are into single forward pass, region proposal network (RPN) is employed for faster process |
| MTCNN | Concentrates on five important landmarks (eyes, nose, mouth) |

Inference Time and Objectness Score are evaluated for different face detection methods using different resolutions for images in the datasets.

TABLE IV
INFERENCE TIME OF TRAINING SET

| Resolutions | Training Datasets Inference Time (s) of Various CNN Methods | | |
|---|---|---|---|
| | YOLOv8 | FASTER R-CNN | MTCNN |
| 40 * 40 | 13.7435 | 179.2755 | 115.4528 |
| 60 * 60 | 9.2719 | 192.4697 | 108.4173 |
| 160 * 160 | 12.3560 | 181.688 | 174.4039 |
| 224 * 224 | 47.584 | 177.551 | 169.6844 |

TABLE V
INFERENCE TIME OF TEST SET

| Resolutions | Testing Datasets Inference Time (s) of Various CNN Methods | | |
|---|---|---|---|
| | YOLOv8 | FASTER R-CNN | MTCNN |
| 40 * 40 | 62.1817 | 988.0068 | 18.8029 |
| 60 * 60 | 53.6966 | 979.5179 | 28.19953 |
| 160 * 160 | 345.599 | 1032.3828 | 32.9586 |
| 224 * 224 | 536.555 | 943.7335 | 36.0668 |

TABLE VI
OBJECTNESS SCORE OF TRAINING SET

| Resolutions | Training Datasets Objectness Score (%) of Various CNN Methods | | |
|---|---|---|---|
| | YOLOv8 | FASTER R-CNN | MTCNN |
| 40 * 40 | 79.9757 | 98.1973 | 99.1803 |
| 60 * 60 | 82.004 | 98.7429 | 99.5659 |
| 160 * 160 | 91.535 | 98.8786 | 99.6416 |
| 224 * 224 | 97.2143 | 98.0715 | 99.8001 |

TABLE VII
OBJECTNESS SCORE OF TEST SET

| Resolutions | Testing Datasets Objectness Score (%) of Various CNN Methods | | |
|---|---|---|---|
| | YOLOv8 | FASTER R-CNN | MTCNN |
| 40 * 40 | 36.7167 | 97.357 | 0 |
| 60 * 60 | 61.91021 | 97.9028 | 85.4735 |
| 160 * 160 | 48.6771 | 97.2944 | 95.6906 |
| 224 * 224 | 44.9968 | 97.1953 | 96.4891 |

It is evident from Table IV, V, VI, and VII, that YOLOv8 achieves low confidence at low resolutions and confidence improves at high resolutions. The Faster R-CNN method provides better confidence at all resolutions with the high computational (or inference) time. On contrast to this, MTCNN achieves the best results with landmark points even at above 60*60 resolutions with tolerable computation time.

*C. Time Measures Visualization*

Inference measures are calculated based on comparison of face detection algorithms. The datasets are processed with resizing as per the required resolution. Once datasets are processed, the beginning and ending of the training network is calculated to know the computation time taken by the model for face detection.
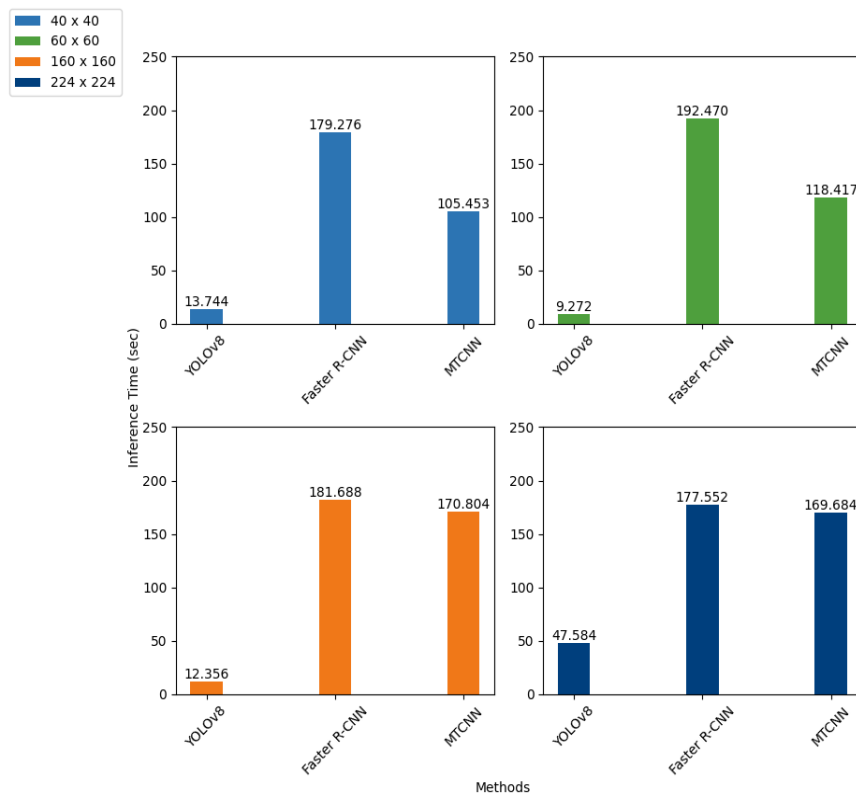
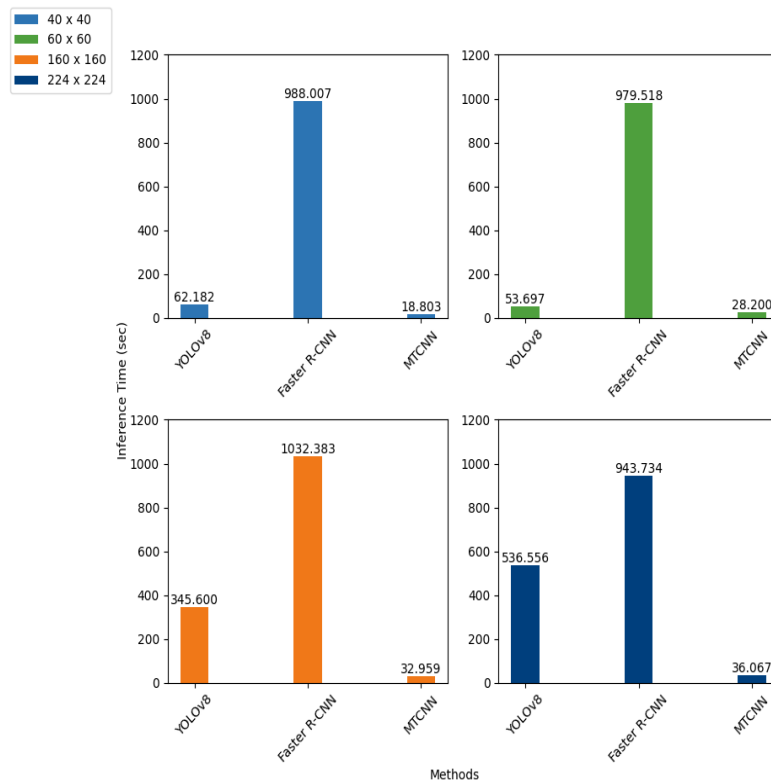Fig.8    Variation of Time measures for training datasets using various CNN Methods



Fig.9    Variation of Time measures for test datasets using various CNN methods

From Fig 8 and Fig 9, it is inferred that the time measure visualization of training and test datasets processed by MTCNN outperforms the other pre-trained algorithms with moderate computational time at any resolution despite its computational complexity.

*D. Confidence Visualization*

The computation of confidence measures incorporates different facial recognition systems. Based on the training network output, fitness or confidence score is evaluated to know the model's average precision of face detection based on the desired formula.
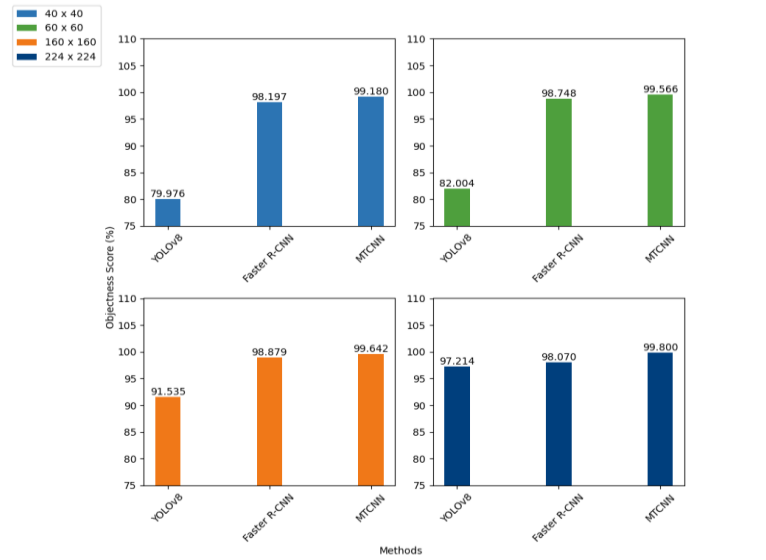


Fig 10. Variation of Score measures for training datasets using various CNN methods
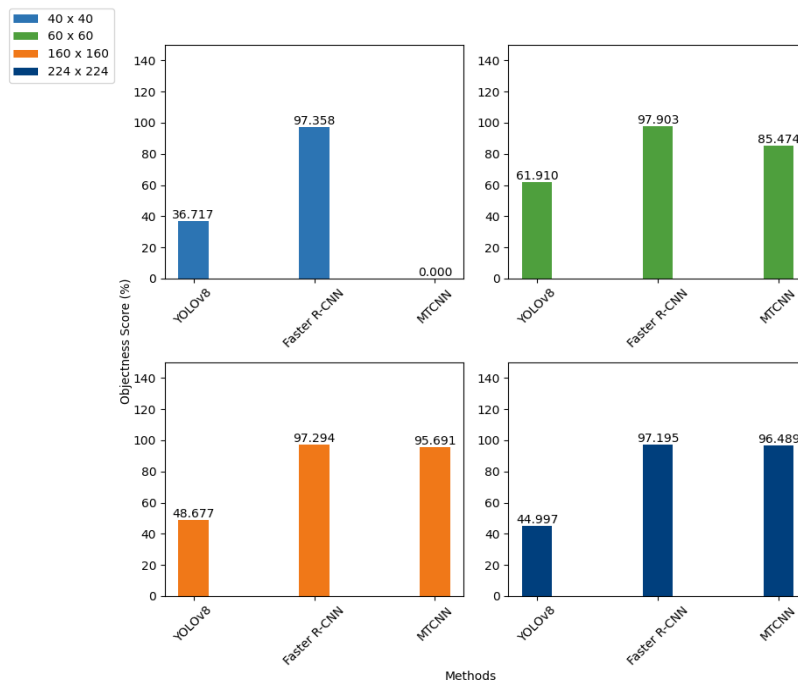


Fig 11. Variation of Score measures for test datasets using various CNN methods

It is observed from Fig 10 and 11, that the score measure visualization of training and test datasets processed by MTCNN performs better than the other pretrained algorithms with better score at low resolution of 60*60 pixels at moderate computation time.

*E. Real-Time Analysis for Face Detection*

Real Time image datasets with confidence scores are analysed by various face detection algorithms. Some of the real time images with different angles and occlusions are experimented by various detection techniques and the results showcase the detected facial images with the boundary boxes and confidence value.

For the real time faces images of various occlusions like angle variation, less lightening effect, specs, the following results have been captured and depicted by the boundary boxes highlighted with corresponding confidence values after the application of IoU and Non max suppression.

*1) Multitask Cascaded Convolutional Neural Networks*:



Fig 12. MTCNN Based Real Time Face Detection at various scenes

Fig 12 depicts that the MTCNN algorithm works better and detects the faces with better background subtraction in the real-time. From the above Fig 12, it is evident that image 2 has a frontal face angle with confidence of 99.4%. In Image 1, multiple face objects are detected by MTCNN with high confidence values of 97.2% and 99.9%.
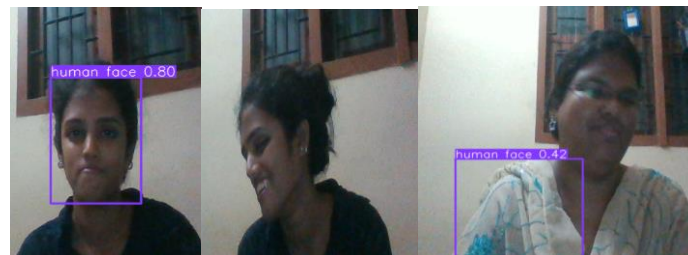
*2) YOLOv8:*



Fig 13. YOLOv8 Based Real Time Face Detection at various scenes

In Fig 13, YOLOv8 algorithm detection results show the possibility of occurrence of face detection with false positive and true negative in image 2 and image 3 respectively. The first image shows a high confidence of 80% and considered above the threshold value. The image 2 from Fig 13 shows that the YOLOv8 missed the detection of facial images at the non-frontal images due to distribution focal loss.

3) *Faster R-CNN:*



Fig 14. Real Time Face Detection using Faster R-CNN at various scenes

From Fig 14, the Faster R-CNN algorithm detection results show that face detection is achieved with better confidence of value 97% irrespective of occlusions with larger inference timing.

TABLE VIII
INFERENCE TIME AND OBJECTNESS SCORE OF UNTRAINED TEST SET

| Methods | Sample Untrained Set Images | |
| | 512 * 512 | |
| | Inference Time (sec) | Objectness Score (%) |
| --- | --- | --- |
| YOLOv8 | 3060.823 | 35.8833 |
| FASTER R-CNN | 1000.3 | 24.3962 |
| MTCNN | 40.0379 | 96.3325 |

Tested face detection models on the sample dataset with trained resolution have achieved high confidence.
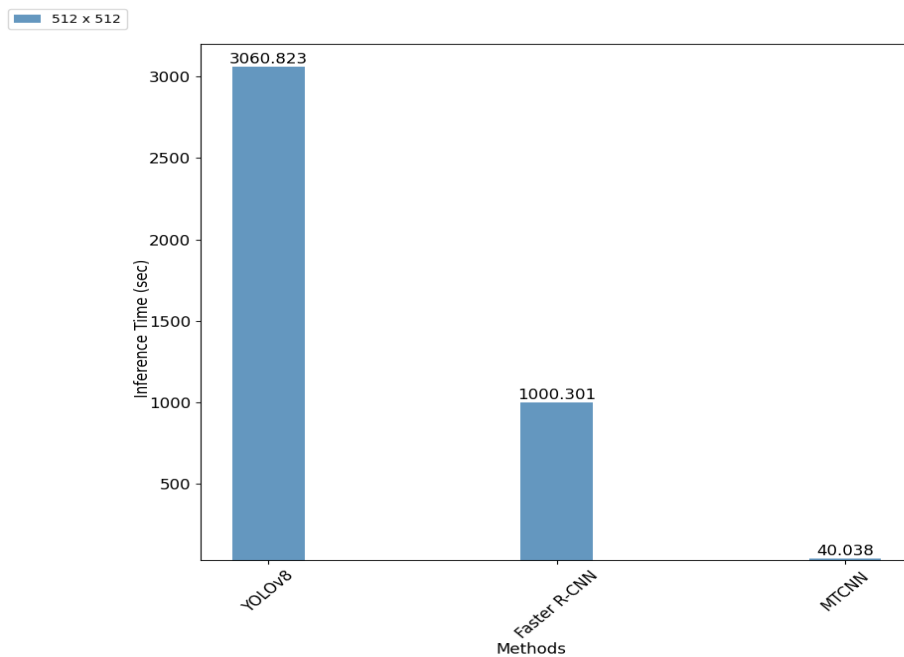


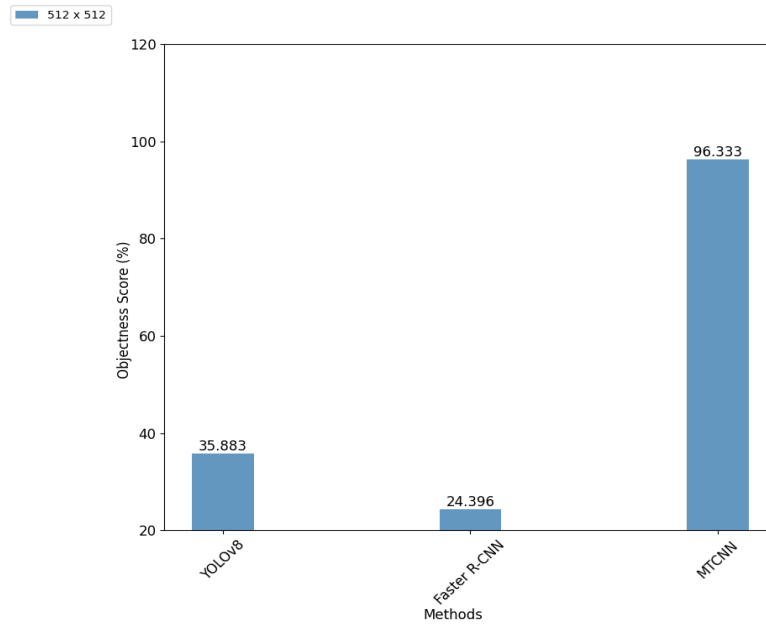Fig 15. CNN methods based Time measures on untrained resolution of images

Fig 16. CNN methods based Score measures on untrained resolution of images

An attempt to test the sample data with untrained resolution of 512*512 is carried out and the behaviour of the models are captured and tabulated in Table VIII.

From Fig 15 and Fig 16, it is observed that even on the untrained resolution of facial image detection, MTCNN performs better than the other two algorithms and takes moderate computation time.

## 5. ETHICAL CHALLENGES IN FACE DETECTION FOR PSYCHIATRY

To address the ethical issues involving privacy, consent, and bias in face detection models, particularly in sensitive applications like psychiatry, and to improve the practical value of proposed methodologies, its compromising to take a comprehensive, multi-faceted approach. The aim is to ensure the models are fair, transparent, secure, and beneficial for all users, while minimizing risks related to misdiagnosis, bias, and invasion of privacy. There are many ways to achieve it:

1. Addressing Bias and Enhancing Fairness in Face Detection Models

    To address bias and enhance fairness in models, some of the strategies to include are having diverse training datasets, using techniques like re-weighting the training data, applying adversarial training, or employing fairness constraints during model optimization to reduce biased outcomes to ensure model working across diverse demographic groups.

2. Improving transparency and interpretability

    To enhance the trustworthiness and practical value of face detection models, especially in sensitive settings like psychiatry, transparency and interpretability are key. This allows both patients and healthcare professionals to understand how the model works and what decisions are being made.

3.  Ensuring Privacy and Secure Data Practices

    Given the sensitive nature of psychiatric data, privacy and data security must be paramount. To safeguard personal information and maintain ethical standards in data handling, data minimization, implement strict access controls to ensure that only authorized individuals have access to data.

4.  Improving Practical Utility Through Context-Awareness

    Face detection systems in psychiatry need to be context-aware to be truly valuable in practice. The same facial expression might have different meanings depending on the context such as whether a patient is experiencing a mental health crisis or is in a therapeutic session.

5.  Promoting Ethical Governance and Oversight

    Finally, to ensure that face detection systems are being used responsibly, it's essential to establish strong governance frameworks.

By integrating fairness, transparency, privacy protections, contextual awareness, and ethical oversight, face detection models can be both more balanced and more practically useful in sensitive applications like psychiatry. Such an approach ensures that these technologies enhance rather than undermine the quality of care, while respecting patient autonomy and minimizing potential harm.

## 6. CONCLUSION

This project has proposed different methodologies for finding the best working model and analysing the efficiencies on accuracy, loss and computational time. In this work, various experiments are performed on the different face detection methods like Faster R-CNN, YOLOv8 and MTCNN neural network based algorithms.

All the face detection methods have been tested on four different resolutions of images. The algorithms are also tested on untrained resolution (512 *512) to test the performance. From the experimental test results,

1.  Faster R-CNN provided consistent accuracy, but as two stage object detection, its high complex computations and iterative RPNs, it consumes double the time of MTCNN when resolution gets higher.
2.  YOLOv8 proves to achieve results at an efficient time than any other methods, but as an object detection algorithm it fails to achieve less true negatives. For higher resolutions without model training, it takes more time also, accuracy gets lower.
3.  MTCNN proved higher accuracy at all resolutions with tolerable inference time comparatively, but fails to detect below 60*60 resolution images.

Based on the real-time face datasets, pretrained detectors like MTCNN, YOLOv8, Faster R-CNN are tested. The observed conclusions are

1.  The MTCNN works for better facial detection of utmost low resolution of 60*60 pixels but fails at below 60*60 resolutions.
2.  The MTCNN results with moderate computation time for processing but has complex computational networks which increases execution load.
3.  Despite the accuracy and computational time, the network is complex and recognition needs additional neural networks to extract and use the feature vectors to implement the recognition task of facial detection.

From the conclusion, it is understood, in accordance with the consistent accuracy, there are still open challenges that exist in face detection. Also, deep learning-based face detection has significantly improved the accuracy, speed, and robustness of face detection systems, making it a valuable tool in diverse applications such as security, healthcare, and human-computer interaction.

In the context of psychiatry, these advancements offer great potential for non-invasive patient monitoring, emotion recognition, and supporting therapeutic interventions. Face detection systems, when integrated with other psychological analysis tools, could greatly enhance diagnostic and treatment processes in mental health settings. Additionally, these technologies offer opportunities for improving the understanding of emotional states, social behaviour, and cognitive functions, which are crucial areas in psychiatric research. However, challenges such as dataset biases, real-time performance issues, and ethical concerns surrounding privacy and consent must be addressed to fully harness the potential of these systems in clinical and research settings.

## 7. FUTURE SCOPE

With the help of face detection information, various neural network classifiers will be analysed for better performance to classify the emotions with the detected face region. The following will be targeted for the future scope in testing classifiers

- Classification of various emotions based on the presence of face regions.
- Provision of high confidence score even at low resolutions.
- Maintenance of tolerable inference time with minimum loss while training the model.

Among the face object detection algorithms, MTCNN works better with face alignment and facial landmarks output. It helps with identifying the face features from the input images. For recognition of facial emotions work, along with the pre-training MTCNN detector, neural network models with multilayers of neuron architecture, will be explored for classification to analyse face emotions by addressing challenges and experiment with those FER systems in real-time. Also, Future work should focus on refining algorithms, improving data diversity, and ensuring the ethical application of these technologies in sensitive environments.

In summary, deep learning-based face detection has made remarkable strides in the past decade, and its continued development holds promise for enhancing psychiatric research and clinical practice. By bridging the gap between technology and psychiatry, these systems could offer valuable tools for clinicians and researchers in assessing and supporting mental health patients.

### REFERENCES

[1]   Xu, Xinzheng & Du, Meng & Guo, Huanxiu & Chang, Jianying & Zhao, Xiaoyang. (2021). *Lightweight FaceNet Based on MobileNet*. International Journal of Intelligence Science. 11. 1-16. 10.4236/ijis.2021.111001.

[2]   F. Zhang, Q. Li, Y. Ren, H. Xu, Y. Song and S. Liu, "*An Expression Recognition Method on Robots Based on Mobilenet V2-SSD*," 2019 6th International Conference on Systems and Informatics (ICSAI), Shanghai, China, 2019, pp. 118-122, doi: 10.1109/ICSAI48974.2019.9010173.

[3]   Nan, Y., Ju, J., Hua, Q., Zhang, H., & Wang, B. (2021). A-MobileNet: *An approach of facial expression recognition*. Alexandria Engineering Journal.

[4]   Szegedy, Christian & Ioffe, Sergey & Vanhoucke, Vincent & Alemi, Alexander. (2016). *Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning*. AAAI Conference on Artificial Intelligence. 31. 10.1609/aaai.v31i1.11231.

[5]   Howard, Andrew & Zhu, Menglong & Chen, Bo & Kalenichenko, Dmitry & Wang, Weijun & Weyand, Tobias & Andreetto, Marco & Adam, Hartwig. (2017). *MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications*.

[6]   Islam, Towfiqul & Ahmed, Tanzim & Rashid, A. & Islam, Taminul & Rahman, Md & Habib, Md. (2022). *Convolutional Neural Network Based Partial Face Detection*. 1-6. 10.1109/I2CT54291.2022.9825259.

[7]   Prameela Naga, Swamy Das Marri, Raiza Borreo, *Facial emotion recognition methods, datasets and technologies: A literature survey, Materials Today*: Proceedings, Volume 80, Part 3, 2023, Pages 2824-2828, ISSN 2214-7853, https://doi.org/10.1016/j.matpr.2021.07.046.

[8]   Szegedy, C., et al. (2015) *Going Deeper with Convolutions*. 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Boston, MA, 7-12 June 2015, 1-9. https://doi.org/10.1109/CVPR.2015.7298594

[9]   Jonathan Huang, Vivek Rathod, Chen Sun, Menglong Zhu, Anoop Korattikara, Alireza Fathi, Ian Fischer, Zbigniew Wojna, Yang Song, Sergio Guadarrama, et al. *Speed/accuracy trade-offs for modern convolutional object detectors.* In CVPR, 2017

[10]  Komar, Myroslav & Yakobchuk, Pavlo & Golovko, Vladimir & Dorosh, Vitaliy & Sachenko, Anatoliy. (2018). *Deep Neural Network for Image Recognition Based on the Caffe Framework*. 102-106. 10.1109/DSMP.2018.8478621.

[11]  Zhang, K., Zhang, Z., Li, Z., & Qiao, Y. (2016). *Joint face detection and alignment using multitask cascaded convolutional networks*. IEEE Signal Processing Letters, 23(10), 1499-1503.

[12]  Lal, Madan & Kumar, Kamlesh & Hussain, Rafaqat & Maitlo, Abdullah & Ruk, Sadaquat & Shaikh, Hidayatullah. (2018). *Study of Face Recognition Techniques: A Survey*. International Journal of Advanced Computer Science and Applications. 9. 10.14569/IJACSA.2018.090606.

[13]  Li, H., Lin, Z., Shen, X., Brandt, J., & Hua, G. (2015). *A convolutional neural network cascade for face detection*. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 5325-5334)

[14]  F. Schroff, D. Kalenichenko and J. Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering," 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 815-823, 2015.

[15]  Jiang, B., Ren, Q., Dai, F., Xiong, J., Yang, J., Gui, G. (2020). *Multi-task Cascaded Convolutional Neural Networks for Real-Time Dynamic Face Recognition Method*. In: Liang, Q., Liu, X., Na, Z., Wang, W., Mu, J., Zhang, B. (eds) Communications, Signal Processing, and Systems. CSPS 2018. Lecture Notes in Electrical Engineering, vol 517. Springer, Singapore. https://doi.org/10.1007/978-981-13-6508-9_8.

[16]  Aaronson, R. Y., Chen, W., & Benuwa, B. B. (2017). *Robust Face Detection using Convolutional Neural Network*. International Journal of Computer Applications, 170(6), 14-20.

[17]  T. Smith, J. Doe, and K. Brown, "Real-Time Face Detection and Tracking with Deep Learning," *IEEE Transactions on Multimedia*, vol. 24, pp. 168-182, 2022.

[18]  L. Zhang, Y. Li, and J. Zhao, "An Improved Face Recognition Method Based on Hybrid Features," *IEEE Access*, vol. 10, pp. 2345-2356, 2022.

[19]  R. Patel and M. Shah, "A Novel Framework for Face Recognition Using GANs," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 9, pp. 4123-4135, Sep. 2022.

[20]  H. Wang, X. Liu, and Y. Zhang, "Face Detection and Recognition in Videos Using Deep Learning," *IEEE Transactions on Image Processing*, vol. 30, pp. 4231-4242, 2022.

[21]  D. Singh and A. Kumar, "Face Recognition Using Ensemble Learning Techniques," *IEEE Transactions on Information Forensics and Security*, vol. 17, no. 1, pp. 112-125, Jan. 2022.

[22]  Q. Xu et al., "Adversarial Attacks on Face Recognition Systems: A Survey," *IEEE Access*, vol. 10, pp. 4567-4580, 2022.

[23]  Y. Zhang, S. Wang, and L. Zhang, "Cross-Age Face Recognition Using Deep Learning," *IEEE Transactions on Image Processing*, vol. 31, pp. 5432-5443, 2022.

[24]  A. R. Khan, P. P. Singh, and R. C. Gupta, "Face Recognition in Unconstrained Environments: A Survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence,* vol. 44, no. 5, pp. 345-367, May 2022.

[25]  T. Liu et al., "Facial Image De-noising for Face Recognition Enhancement," *IEEE Transactions on Information Forensics and Security,* vol. 17, no. 3, pp. 654-667, Mar. 2022.

[26]  M. K. Chen and J. F. Liu, "Face Recognition Based on Hybrid Deep Learning Models," *IEEE Transactions on Systems, Man, and Cybernetics: Systems,* vol. 52, no. 2, pp. 789-802, Feb. 2022.

[27] X. Liu and D. Zhang, "Dynamic Face Recognition in Surveillance Videos," *IEEE Transactions on Image Processing,* vol. 31, pp. 876-887, 2022.

[28] Y. Guo et al., "High-Performance Face Recognition on Mobile Devices," *IEEE Access*, vol. 10, pp. 3245-3258, 2022.

[29] R. Sharma and P. Sharma, "Facial Recognition Using 3D Morphable Models," *IEEE Transactions on Image Processing*, vol. 31, pp. 1200-1212, 2022.

[30] J. Smith and A. Brown, "Face Detection Using Dlib for Real-Time Applications," *IEEE Trans. Image Process.*, vol. 29, no. 4, pp. 1234-1242, Apr. 2020.

[31] M. Johnson et al., "Comparative Analysis of MTCNN and Dlib for Face Detection," *IEEE Access*, vol. 8, pp. 5678-5685, Jan. 2021.

[32] L. Zhang and K. Lee, "Improving Face Recognition Accuracy with FaceNet and MTCNN," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 9, pp. 3678-3685, Sep. 2021.

[33] P. Wang, R. Kumar, and S. Gupta, "Real-Time Face Detection Using Dlib and OpenCV," *IEEE Open Journal of Computing. Inteligence.*, vol. 2, pp. 45-52, Jul. 2022.

[34] A. Davis, "FaceNet-Based Face Recognition for Security Systems," *IEEE Trans. Inf. Forensics Security,* vol. 15, pp. 123-134, Feb. 2020.

[35] T. Wilson and C. Smith, "Application of MTCNN for Facial Feature Extraction," *IEEE Trans. Multimedia,* vol. 23, no. 3, pp. 456-464, Mar. 2022.

[36] H. Kim, "Enhancing Face Detection Accuracy with Dlib and Machine Learning Techniques," *IEEE J. Sel. Areas Commun.,* vol. 38, no. 5, pp. 1010-1017, May 2020.

[37] R. Patel and J. Lee, "A Survey of Face Detection Methods: Dlib, MTCNN, and Beyond," *IEEE Access,* vol. 10, pp. 3456-3469, Mar. 2023.