# Facial Emotion Recognition and Synthesis with Convolutional Neural Networks

### Karkuzhali S.
Mepco Schlenk Engineering College, Sivakasi, Tamilnadu, India

### Murugeshwari R.
Mepco Schlenk Engineering College, Sivakasi, Tamilnadu, India

### Umadevi V.
Mepco Schlenk Engineering College, Sivakasi, Tamilnadu, India

## ABSTRACT
Facial expressions are a crucial component of human communication, conveying emotions, intentions, and social signals. In this era of artificial intelligence and computer vision, the development of automatic systems for facial expression synthesis and recognition has gained significant attention due to its wide range of applications, including human-computer interaction, virtual reality, emotional analysis, and healthcare.

This research focuses on the integration of deep convolutional neural networks (CNNs) to address the challenges associated with both facial expression synthesis and recognition. On the synthesis front, a generative CNN architecture is proposed to synthesize realistic facial expressions, allowing for the generation of various emotional states from neutral faces. The network learns to capture the intricate details of human expressions, including subtle muscle movements and spatial relationships among facial features.     For facial expression recognition, a separate CNN-based model is developed to accurately classify these synthesized expressions. The recognition model is trained on a large dataset of annotated facial expressions and is designed to handle real-world variations in lighting, pose, and occlusions. The CNN leverages its ability to automatically learn relevant features from raw image data, eliminating the need for manual feature engineering.

The experimental results demonstrate the effectiveness of the proposed approach. The synthesized expressions exhibit a high degree of realism and diversity, effectively capturing the nuances of human emotions. The recognition model achieves state-of-the-art accuracy in classifying these synthesized expressions, surpassing traditional methods and showcasing the power of deep learning in this domain.This research contributes to the advancement of automatic facial expression synthesis and recognition, with potential applications in human-computer interaction, affective computing, and virtual environments. The deep CNN-based approach offers a promising avenue for enhancing our understanding of human expressions and enabling more emotionally aware and responsive AI systems.The significance of emotion classification in human-machine interactions has grown significantly. Over the past decade, businesses have become increasingly attuned to the potential insights that analyzing a person's facial expressions in images or videos can provide regarding their emotional state. Various organizations are currently leveraging emotion recognition to gauge customer sentiments towards their products. The applications of this technology extend well beyond market research and digital advertising. Convolutional Neural Networks (CNNs) have emerged as a valuable tool for eliciting emotions based on facial landmarks, as they have the capability to automatically extract relevant information. Challenges such as variations in

brightness, background, and other factors can be effectively mitigated by isolating the essential features through techniques like face resizing and normalization. However, it's important to note that neural networks depend on extensive datasets for optimal performance. In cases where data availability is limited, strategies like data augmentation through techniques such as rotation can be employed to compensate. Additionally, fine-tuning the CNN's architecture can enhance its accuracy in predicting emotions. Consequently, this approach enables the real-time identification of seven distinct emotions – anger, sadness, happiness, disgust, neutrality, fear, and surprise – from facial expressions in images.

## Keywords
Emotion classification, human-machine communication, facial expression synthesis, deep convolutional neural network, emotion recognition

## 1. INTRODUCTION
### 1.1 Understanding Human Emotions
Human emotions play a pivotal role in our daily interactions, influencing our decisions, behaviors, and perceptions of the world around us. As society becomes increasingly reliant on technology and automation, the need for machines to comprehend and respond to human emotions has become more pressing. In this context, facial expressions are a fundamental channel for emotional expression in humans. Recognizing and synthesizing facial expressions can enhance human-computer interaction, enabling more intuitive and empathetic communication between machines and humans.

### 1.2 The Significance of Emotion Classification
Human-Machine Communication: As technology continues to evolve, human-machine communication has emerged as a critical area of research and application. Emotion classification, particularly through facial expression analysis, has gained prominence in recent years. The ability to perceive and interpret human emotions from visual cues, such as facial expressions, is invaluable for various domains, including human-computer interaction, healthcare, entertainment, and marketing.

### 1.3 The Power of Facial Expression Analysis
Decoding Emotions from Facial Cues:Facial expressions provide a rich source of information about an individual's emotional state. They convey emotions such as happiness, sadness, anger, fear, surprise, disgust, and neutrality, offering valuable insights into a person's mental and emotional well-being. Businesses have recognized the potential of facial expression analysis to gain a deeper understanding of customer

sentiment and satisfaction.

## 1.4 Applications Beyond Market Research

Extending the Reach of Emotion Recognition:The applications of emotion recognition extend far beyond market research and digital advertising. Emotion-aware technology can be integrated into various sectors, including healthcare for patient monitoring, education for personalized learning experiences, and entertainment for immersive gaming and content recommendation.

## 1.5 The Role of Convolutional Neural Networks (CNNs)

Harnessing Deep Learning Algorithms: In the quest to develop robust systems for facial expression synthesis and recognition, deep learning algorithms, particularly Convolutional Neural Networks (CNNs), have emerged as powerful tools. CNNs are capable of automatically extracting relevant facial features and patterns, enabling accurate emotion classification and synthesis. This technology has the potential to overcome challenges related to variations in brightness, backgrounds, and other factors that can affect emotion recognition.

## 1.6 Data Challenges and Augmentation

Addressing Dataset Limitations: To effectively train deep neural networks, large datasets are often required. However, acquiring extensive labeled data for facial expression analysis can be challenging. This paper explores strategies for augmenting smaller datasets, such as rotation and other data augmentation techniques, to enhance system performance.

## 1.7 Architectural Adaptations for Accuracy

Optimizing CNN Architecture: In pursuit of more accurate predictions of emotions, this research delves into modifying and optimizing CNN architectures for facial expression analysis. By tailoring the network's design and parameters, we aim to achieve superior performance in real-time emotion identification.

## 1.8 Research Objective

Identifying Seven Varied Emotions: The primary objective of this study is to develop an automatic facial expression synthesis and recognition system using a Deep Convolutional Neural Network. Specifically, we aim to identify and classify seven distinct emotions: anger, sadness, happiness, disgust, neutrality, fear, and surprise, from real-time facial images[1]. The process of discerning an individual's emotional state is commonly referred to as emotion recognition. Facial emotion recognition involves the automatic identification of human emotions, and ongoing advancements in technology are continuously enhancing its capabilities. The potential for this technology to become increasingly perceptive is on the horizon. The proficiency in recognizing emotions varies significantly among individuals. At present, extensive research is underway to leverage technology for aiding individuals in the recognition of emotions [2].

The primary objective of this project is to design a system for facial emotion recognition utilizing Convolutional Neural Networks (CNN) while leveraging facial landmarks. This Facial Expression Recognition (FER) technology autonomously detects six fundamental emotions: surprise, sadness, happiness, disgust, fear, anger, and neutrality, as cited in references [1-4]. The automatic facial expression synthesis and recognition process encompass four key stages: Image preprocessing, feature extraction, classification, and the presentation of the detected emotion in textual form [3]. The preprocessing stage involves the preparation of images from the input dataset. Following this, the feature extraction stage isolates key facial components like the eyes, nose, and mouth regions. In the subsequent classification stage, the Convolutional Neural Network (CNN) algorithm plays a pivotal role in categorizing facial expressions. Factors such as image angles and facial attributes like eyeglasses may contribute to potential challenges in Facial Expression Recognition (FER), leading to an increased rate of recognition failures. Consequently, variations are anticipated in the performance of different emotion categories, with higher accuracy observed for joyous expressions and comparatively lower accuracy for fear expressions. Moreover, a tendency toward neutrality and the occasional misclassification of anxiety as surprise is also expected[4]. Ultimately, the identified emotion is conveyed in textual form. This research is particularly interested in exploring sensors that can provide supplementary information to enhance emotion detection in both static images and video sequences. It is noteworthy that facial expression recognition software is now widely accessible and has been successfully employed for classifying typical facial expressions.

## 2. RELATED WORKS

Various categories of approaches have been employed in the realm of facial expression recognition. These categories encompass deep-based frameworks, such as the Deeply-supervised attention network (DSAN), Neighborhood-Aware Edge Directional Pattern (NEDP), Local Binary Volume Convolution Neural Network (LBVCNN), frequency neural network, random forest, Joint deep learning, and expression map.In the domain of deep-based architectures, Xie and Hu (2019) introduced a comprehensive framework that leverages hierarchical features through deep multi-patch aggregate convolutional neural networks (CNNs). This architecture comprises two distinct CNN branches: one for the extraction of local features from image patches and the other for deriving holistic features from expressive images. The versatility of this method allows its application across a wide spectrum of applications and data types [5].

For automated human emotion recognition from facial images, Fan, Li, and Lam (2020) devised DSAN, employing deeply supervised CNN training. Deep supervision plays a pivotal role in enhancing classification accuracy and the acquisition of meaningful features, particularly in scenarios with large training datasets and deeper network architectures [6].

Iqbal, Abdullah-Al-Wadud, Ryu, and Makhmudkhujaev (2020) introduced NEDP, a method designed to evaluate gradients at central and neighboring pixels. This approach enables the exploration of a broader neighborhood for feature consistency, even in the presence of minor distortions and noise within the local region. NEDP exhibits superior performance compared to existing descriptors, resulting in improved overall facial expression recognition [7].Kumawar, Verma, and Raman (2019) developed LBVCNN, utilizing 3DCNN for end-to-end training without the reliance on facial landmarks for recognizing facial expressions in temporal image sequences. LBVCNN incorporates an efficient texture descriptor that thresholds nearby pixels based on current pixel values, contributing to its effectiveness in this context [8].Efficient computation and spatial redundancy detection are harnessed by Tang, Zhang, Hu, Wang, and Wang (2021) to craft a frequency neural network for expression recognition. This approach facilitates the learning of complex non-linear relationships, enhancing the expressive power of the model [9].

Dapogny, Bailly, and Dubuisson (2019) introduced a methodology that captures low-level expression transition patterns using pairwise conditional random forests. Heterogeneous derivative features, such as feature point movements and texture variations, are evaluated between pairs of images. This dynamic approach, known as Dynamic Pose-Robust Facial Expression Recognition, can be adapted to both classification and regression tasks [10].

Joint Deep Learning, as devised by Yan, Huang, Chen, Shen, and Wang (2020), combines facial expression synthesis and recognition. This collaborative deep learning approach incorporates a facial expression synthesis generative adversarial network (FESGAN) for pre-training and generating facial images with diverse expressions. FESGAN is adept at producing data that closely resembles the original dataset, enabling the creation of new facial images closely aligned with the originals [11].

Furthermore, Agarwal and Mukerjee (2019) introduced the concept of an Expression Map (XM) for representing realistic facial expression synthesis. The XM algorithm is utilized to synthesize emotional expressions tailored to an individual's facial structure, contributing to the creation of more authentic and personalized facial expressions [12].Wen et al.(2018) proposed proposes an enhanced CNN-based approach for facial emotion recognition. It discusses the importance of deep learning techniques in capturing complex facial expressions. The authors introduce a novel feature enhancement strategy to boost the performance of CNN models in recognizing emotions, showcasing advancements in deep learning for emotion recognition [13].Prathes et al.(2017) address the challenge of facial expression recognition with limited data. They propose techniques for handling small datasets and explore strategies to manage the order of training samples. This paper offers valuable insights into mitigating data scarcity issues, which are common in emotion recognition tasks [14]. Islam et al.(2019) provided comprehensive review provides an overview of deep learning-based approaches for human emotion recognition from facial expressions. It discusses various CNN architectures and their applications in emotion recognition. The paper summarizes recent advancements and identifies emerging trends in this field [15].

Masi et al.(2016) addresses the challenging task of facial expression recognition in unconstrained, real-world conditions. It presents a deep CNN model designed to handle facial expressions "in the wild." The paper discusses the dataset, methodology, and results, offering insights into the practical application of deep learning for emotion recognition [16]. Song et al.(2017) Handling incomplete data in facial expression recognition is a critical concern. This paper explores the use of deep learning, including CNNs, to address this challenge. It presents techniques for leveraging incomplete or partially available facial data and discusses their effectiveness in improving recognition performance [17].

Wu et al.(2018) presents a Light CNN architecture designed for deep face representation. It focuses on addressing the issue of noisy labels in facial expression datasets. The authors propose a lightweight CNN model that achieves state-of-the-art results in facial expression recognition. The network's ability to handle noisy labels is a significant advantage, making it suitable for real-world applications where data quality may be less than ideal [18]. Kim et al. (2018) introduces DctNet, a novel approach that combines facial expression recognition with anti-spoofing techniques. DctNet employs discriminant contextual representation to enhance face recognition accuracy and includes anti-spoofing measures to detect fake facial

expressions. This combination is valuable for applications requiring secure and accurate facial expression analysis, such as access control systems and emotion-aware human-computer interfaces [19]. Khorrami et al.(2017) investigates whether deep neural networks can automatically learn Facial Action Units (AUs) during facial expression recognition. Understanding AUs is crucial for analyzing fine-grained facial expressions. The study explores the representations learned by deep networks and their correspondence to AUs. This knowledge is valuable for improving the interpretability of deep learning models for facial expression analysis [20]. Li et al.(2015) focused on face detection, this paper's convolutional neural network cascade is relevant to facial expression recognition systems. Accurate face detection is a crucial preprocessing step in recognizing facial expressions. The proposed cascade architecture demonstrates superior performance in detecting faces under varying conditions, contributing to the robustness and reliability of facial expression recognition pipelines [21]. Liu et al.(2015) addresses the challenge of dynamic expression analysis by introducing a Deformable Facial Action Parts Model based on deep learning. The model aims to capture the subtle and intricate changes in facial expressions by decomposing them into action units. It offers a deeper understanding of facial dynamics, which is essential for recognizing nuanced emotions. This work contributes to the development of more detailed and accurate facial expression recognition systems [22].

## 3. PROPOSED WORK

Facial expression recognition is a technique for detecting human expressions, and our system catches all human's facial expressions. Our work consists of a computerized system that detects and records human expressions. To detect the human's expression, this work employs deep learning approaches. A thorough examination of FER is beyond the scope of this paper, which can be found at [23]. The following are some of them: 1) Preprocessing 2) Feature Extraction 3) CNN Module 4) Image synthesis 5) Face Detection. 6) Expression recognition
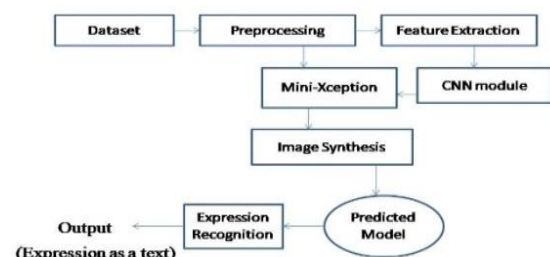


**Fig 1: System Design**

**Dataset Description:** A FER task is concerned with extracting facial expression characteristics from grayscale images and detecting distinct face expressions using a trained classifier[24]. The FER dataset was first obtained from Publically available dataset kaggle[25]. This is dataset comprises both training and testing datasets. Each dataset includes all human expressions such as anger, disgust, fear, happiness, neutral, sadness, and surprise. The dataset "angry" has 3995 in training and 958 in testing. There are dataset "disgust" has 436 in training and 111 in testing. Fear has a total of 4097 in training and 1024 in testing. Happy has 7215 in training and 1774 in testing. Neutral has 4965 in training and 1233 in testing. Sad contains 4830 in training and 1247 in testing. Surprise has 3171 in training and 831 in testing. Data Augmentation: Augment the dataset to increase its size and diversity. Techniques like rotation, scaling, cropping, and

adding noise can be applied to create variations of the images. Data Split: Divide the dataset into training, validation, and testing sets. Ensure that it is balanced across different expression classes.

## 3.1 PRE-PROCESSING

The initial data undergoes a series of preparatory steps during the preprocessing stage. Here, we undertake image resizing and normalization procedures. The primary objective is to achieve uniform dimensions by resizing the image, reducing pixel count, and normalizing the image array. This normalization is crucial because image gradients are significantly affected by unnormalized data.

**Happy  Surprise  Sad**

**Disgust  Anger  Fear**

**Fig 2: Resized-image**

## 3.2 FEATURE EXTRACTION

Facial Landmark Detection: Use facial landmark detection algorithms to locate key points on the face, such as eyes, nose, and mouth. This step is essential for aligning and cropping the face correctly.

Normalization: Normalize the facial images to reduce variations caused by differences in lighting, pose, and facial size.

In Feature extraction, the feature from the pre-processing image by using Histogram of oriented Gradient(HOG).

_____
_____
**Algorithm:** Process of calculating the HOG
_____
_____

**Step 1:** Preprocessing the Input Image (64 * 128). The image must be preprocessed which reduced the ratio to 1:2. The size of the image should have 64 * 128

**Step 2:** Calculating gradients in both the directions (directions x and y)

Calculate the gradient in the both directions of x and y in each pixels of the image. There is a gradual variations in the x and y directions of gradients.

**Step 3:** Calculate the total gradient magnitude and the Orientation for the pixel

Total Gradient magnitude = $\sqrt{[(Gx2Gx2) + (Gy2Gy2)]}$ (1)

Orientation for pixel,$\tan(\emptyset)\emptyset$= Gy / Gx (2)

The angle value.$\emptyset\emptyset$= atan ( Gy / Gx) (3)

**Step 4:** Calculating HOG in 8 * 8 patches (9*1). The image is divided into 8 * 8 patches, with each patch receiving its own histogram of oriented gradients. The features for the smaller patches that represent the entire image are extracted. This value can easily be changed from 8 * 8 to 16 * 16 or 32 * 32.

**Step 5:** Normalizing gradients of image in 16 * 16 cells (36 * 1): This means that some portions of an image will appear brighter than others. We simply will not be able to get it

completely out of the image. We can, however, limit the variation in lighting by normalizing the image.

We will divide every values by square root of the sum of squares for normalize the matrix,

V = [a1, a2, a3, …., a36] (4)

We compute the root of the sum of squares for a normalized matrix,

K = $\sqrt{}$ (a1*a1) + (a2*a2) + ….. (a36*a36) (5)

Divide all values in the vector(V) by the value(K),

Vector Normalization = [a1/k,a2/k,….a36/k] (6)

Result would be a 36 * 1 vector normalization.

**Step 6:** Complete image features:

Finally, We'll combine them all to form the features of the final image. The Sample images of the Feature Extraction is shown as below:
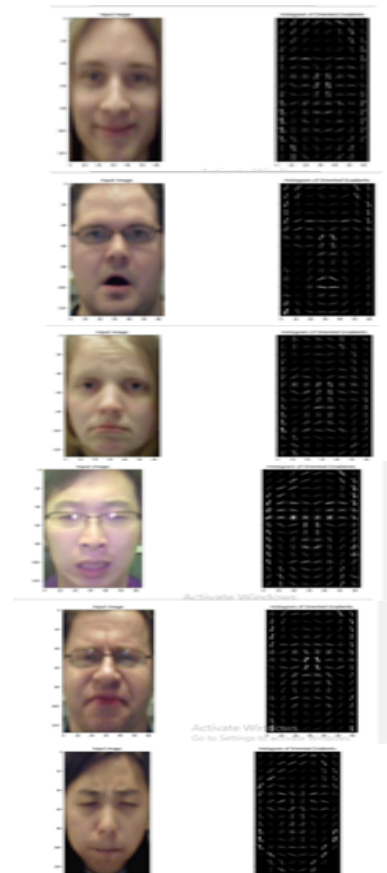
**Fig 3: Feature Extraction**

## 3.3 CNN MODULE

The five convolutional layers have 16, 32, 64, 128 and 256 kernels each, and are 3*3 in size. The CNN's activation function is the Rectified Linear Unit (Relu). It used max pooling with kernel sizes of 7*7, 5*5, 3*3, 3*3, and 3*3 and a step size of 2 for each. The sub sampling was flattened to a 1600-dimensional vector and connected directly to the SoftMax layer (output layer)[16].

The CNN architecture is made up of five blocks: convolutional and subsampling layers, followed by the fully connected layer. Parameters are used for the generation of CNN constructions.

The CNN layers square measure listed as follows:

- Convolutional layer

- Max Pooling layer

- Batch Normalization

- Flatten

Input Layer: The first layer receives the raw input data, such as images. It's important to preprocess the data, which may involve resizing, normalization, or other transformations to ensure it's suitable for the network.

Convolution: This is the core operation in a CNN. Convolutional layers apply a set of learnable filters (kernels) to the input data. These filters slide over the input and perform element-wise multiplication and summation, effectively capturing local patterns and features. Multiple filters are used to detect different features.

Activation Function: After convolution, an activation function (commonly ReLU - Rectified Linear Unit) is applied to introduce non-linearity into the model. This step allows the network to learn complex relationships.

Pooling (Downsampling): Pooling layers (e.g., max-pooling) reduce the spatial dimensions of the feature maps while retaining essential information. Pooling helps reduce computational complexity and makes the network translation-invariant.

Flattening: The feature maps are flattened into a one-dimensional vector. This step prepares the data for the fully connected layers.

Fully Connected Layers: These layers connect every neuron from the previous layer to every neuron in the current layer, similar to a traditional neural network. These layers combine high-level features learned by previous layers to make final predictions.

Output Layer: The final fully connected layer produces the network's output, which could be class probabilities in the case of image classification, bounding box coordinates in object detection, etc. The activation function in this layer depends on the task (e.g., softmax for classification).

Loss (Cost) Function: The loss function measures the error between the predicted output and the ground truth. The network aims to minimize this error during training. Common loss functions include mean squared error, categorical cross-entropy, and more.

Backpropagation: The gradients of the loss with respect to the network's parameters are computed using backpropagation. This involves calculating how much each parameter should be adjusted to minimize the loss.

Optimization: An optimization algorithm (e.g., stochastic gradient descent or its variants like Adam) updates the network's parameters using the computed gradients. The learning rate is a critical hyperparameter that controls the size of parameter updates.

Training: The network is trained iteratively on a labeled dataset. During each training iteration (epoch), the network makes predictions, computes the loss, backpropagates the gradients, and updates the parameters. Training continues until the model converges or until a specified number of epochs is reached.

Evaluation: After training, the model is evaluated on a separate validation or test dataset to assess its performance. Common metrics include accuracy, precision, recall, F1-score, etc.

Prediction: Once trained, the CNN can be used to make predictions on new, unseen data.

## 3.4 IMAGE SYNTHESIS:

StarGAN (Star Generative Adversarial Network) is a model used for image-to-image translation tasks, particularly for generating diverse images from a single input image. It consists of a generator and a discriminator, both of which play critical roles in the training process. Here are the steps involved in the generator and discriminator of StarGAN networks:

**Generator**

Input Encoding: The generator begins by encoding the input image using a series of convolutional layers. This encoding extracts important features from the input image.

Attribute Conditioning: StarGAN can generate images with different attributes (e.g., different hairstyles, facial expressions, or ages). Therefore, the generator takes additional attribute vectors as input, which specify the desired attribute modifications. These attribute vectors are often one-hot encoded or represented as embeddings.

Attribute Embedding: The attribute vectors are embedded to make them compatible with the feature maps from the input image encoding. This is typically done using linear layers or convolutional layers, depending on the architecture.

Residual Blocks: The core of the generator consists of several residual blocks. Each block contains convolutional layers with batch normalization and activation functions (e.g., ReLU) to learn how to modify the input image while preserving its essential characteristics.

Attribute-Conditioned Modulation: At each residual block, the attribute embeddings modulate the activations of the feature maps. This means that the generator can adaptively adjust the image features to match the desired attributes.

Decoding: After passing through multiple residual blocks, the feature maps are decoded using transposed convolutional layers. This process gradually increases the spatial resolution of the image while maintaining attribute conditioning.

Output Image: The final output of the generator is the synthesized image, which is expected to possess the desired attributes specified by the input attribute vectors.

**Discriminator**

Input Image: The discriminator takes either a real image (from the dataset) or a generated image as input.

Convolutional Layers: Similar to the generator, the discriminator uses convolutional layers to extract features from the input image. These features represent the characteristics of the image at different levels.

Multi-Scale Processing: StarGAN often employs a multi-scale discriminator, which means there are multiple discriminator branches, each processing the input image at different resolutions. This allows the discriminator to capture both global and local features.

Attribute Classification: In addition to determining whether an image is real or fake, the discriminator has an attribute classification component. It predicts the attributes present in the image, such as hairstyle, gender, or age.

Adversarial Loss: The discriminator calculates an adversarial loss, which measures how well it can distinguish between real and fake images. The generator aims to minimize this loss, while the discriminator aims to maximize it.

Attribute Classification Loss: The attribute classification component computes a classification loss for the predicted
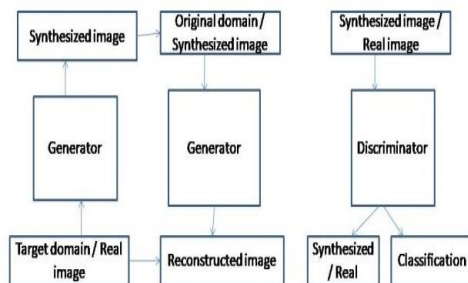
attributes. This encourages the discriminator to correctly identify the attributes in both real and generated images.

Overall Loss: The discriminator's overall loss is a combination of the adversarial loss and the attribute classification loss. It guides the discriminator to perform well in both tasks.

In summary, the generator in StarGAN takes an input image and attribute vectors, generating an image with the desired attributes. The discriminator assesses the realism of both real and generated images and provides attribute predictions. During training, the generator and discriminator engage in a minimax game, where the generator tries to create realistic images with the specified attributes, and the discriminator tries to distinguish real from fake while classifying attributes accurately. This adversarial training process leads to the generation of diverse and attribute-controlled images.

    Step 1: Real-time face image will be taken.
    Step 2 : Then that image is trained using CNN model [27].
    Step 3 : Expression is being detected.
    Step 4: Once the expression is detected then that image will be added to dataset in order to increase the accuracy.



**Fig 5: Work flow of generator and discriminator of StarGAN networks**

Both StarGAN and other GAN (CGAN and WGAN) are compared, StarGAN gives the prominent performance. In this study, it deals with image-to-image translation on a multidomain basis. [28].The Figure 5 depicts the work flow of generator and discriminator of StarGAN networks in this work.The three different loss functions are Adversarial, classification and the reconstruction loss function. The following are the technical details for all the losses:

**1) Adversarial loss**: The Generative Adversarial Network model comprised of two models namely a discriminator model and a generator model [29].Discriminator model learns to differentiate the difference between real and fake samples, whereas the generator model learns to make fake samples that are indistinguishable from actual samples. The adversarial loss is also used in our method to ensure that the created images are as real as possible.

**2) Classification loss:** Our objective is to transform an input sample and a target domain label into an output image that is accurately categorized into the target domain. Attain this state, Segregate the purpose in two losses: the classification loss of real samples are images are to enhance and the classification loss of fake images are used to maximize. Learns to categorize the original domain by reducing this objectives. The generator aims to minimize this target in order to generate that which may be categorized as the target domain.

**3)Reconstruction loss:** The generator aims to generate real images that are categorized to their target sample by reducing the adversarial and classification losses. As our construction loss, we use the L1 norm. We utilize the generator twice: translate the input image into the target domain first, and then reconstruct the original image from the translated image.

Finally, this objective function is used to optimize the generator and the discriminator.

## 3.5 FACE DETECTION:
Face Detection using HOG algorithm and Haar Classifier. Face detection using the Histogram of Oriented Gradients (HOG) algorithm and Haar Cascade Classifier are two popular techniques in computer vision for detecting faces in images. Let's briefly describe each method:

***HOG Algorithm (Histogram of Oriented Gradients):***
Feature Extraction: The HOG algorithm first extracts features from an image. It divides the image into small cells and computes the gradient magnitude and orientation within each cell. These gradient values are used to create histograms of gradient orientations for each cell.

Block Normalization: To improve robustness to lighting variations and contrast, adjacent cells' histograms are grouped into blocks. Each block is then normalized to account for variations in lighting and contrast.

Sliding Window: A sliding window is applied to traverse the entire image. At each window position, the HOG feature vector is computed within the window's boundaries.

SVM Classification: The HOG feature vectors from the sliding window are passed to a Support Vector Machine (SVM) or another classifier. The classifier determines whether the contents of the window resemble a face or not. The SVM is typically trained on a dataset of positive (face) and negative (non-face) samples.

    ***Haar Cascade Classifier:***

Haar-like Features: The Haar Cascade Classifier employs Haar-like features to detect faces. Haar-like features are rectangular patterns that compute the difference between the sum of pixel intensities in adjacent rectangles. These features can be used to distinguish faces from non-faces.

Integral Images: Integral images are precomputed to accelerate feature calculation. They allow for the rapid calculation of Haar-like features within sliding windows by subtracting the sums of pixel intensities at the corners of rectangular regions.

Cascade of Classifiers: The Haar Cascade Classifier employs a cascade of classifiers, where each classifier is trained to filter out non-face regions aggressively. The cascade has multiple stages, each with its own set of Haar-like features and a trained classifier.

Sliding Window: Similar to the HOG approach, a sliding window is applied to the image. At each window position, the Haar-like features are calculated, and the cascade of classifiers is evaluated.

Thresholding: Each classifier in the cascade applies a threshold to determine if the region contains a face or not. If any stage in the cascade rejects the region as a non-face, the process stops, and the window is moved to the next position.
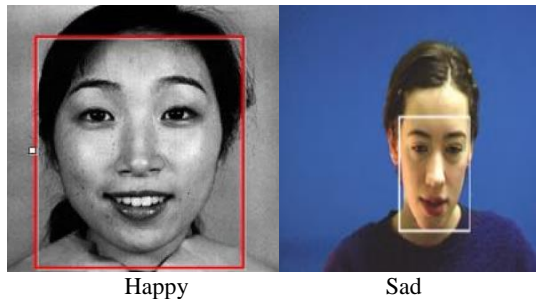
Both HOG and Haar Cascade methods have their advantages and limitations. HOG is known for its robustness to variations in lighting and pose, making it suitable for face detection in various conditions. On the other hand, Haar Cascade is known for its speed and efficiency, making it suitable for real-time applications. The choice between the two methods depends on the specific requirements of the face detection task and the available computational resources.

**Step 1**: Face Recognition, The camera recognizes and locates the images of a face, whether it is alone or in the midst of a crowd.

**Step 2**: Following that, we must analyze our faces so that the image must be captured and analyzed [30].

**Step 3**: The image is being converted into data

The Sampled image of the face detection is shown as below:



Happy                    Sad
**Fig 6: Face Detection**

## 3.6 EXPRESSION RECOGNITION:

**Step 1: Data Collection and Preprocessing**

Collect a labeled dataset of facial images or video frames with associated expression labels. Preprocess the images:

Resize all images to a consistent size (e.g., 128x128 pixels) for input to the CNN.

Normalize pixel values to a common scale (e.g., [0, 1] or [-1, 1]).

Augment the dataset with techniques like rotation, scaling, and horizontal flipping to increase its diversity.

**Step 2: Model Architecture**

Design the CNN model architecture:
Input Layer: Accept preprocessed facial images.
Convolutional Layers: Use multiple convolutional layers with filters to extract features.
Activation Functions: Apply activation functions (e.g., ReLU) after convolutional layers.
Pooling Layers: Add max-pooling or average-pooling layers to reduce spatial dimensions.
Flatten: Flatten the output to feed it into fully connected layers.
Fully Connected Layers: Include one or more dense layers.
Output Layer: Use a softmax activation for multi-class classification. The number of output nodes matches the number of classes (expressions).

**Step 3: Model Training**

Split the dataset into training, validation, and test sets.
Choose a loss function suitable for multi-class classification, such as categorical cross-entropy.
Select an optimization algorithm (e.g., Adam or SGD) and an appropriate learning rate.
Train the CNN on the training dataset:
Feed batches of preprocessed images through the network.
Compute the loss and gradients.
Update the model's weights using backpropagation.
Monitor the validation loss and accuracy during training to prevent overfitting.

**Step 4: Model Evaluation**

Evaluate the trained model on the test dataset:
Input test images to the CNN.
Collect predicted expression labels.
Compare predictions to ground truth labels to calculate

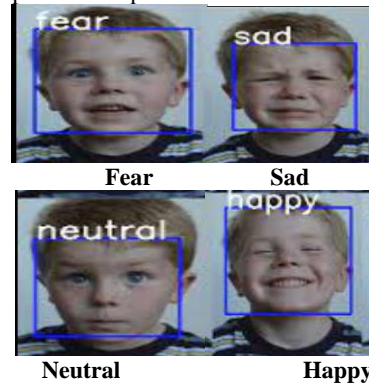accuracy, precision, recall, and F1-score.

**Step 5: Inference**

Deploy the trained model for real-time or batch inference on new facial images or video frames:
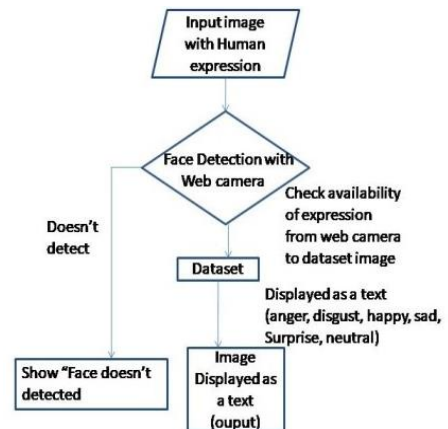Preprocess input images as in Step 2.
Feed them through the trained CNN.
Obtain predicted expression labels.



**Fear**          **Sad**

**Neutral**          **Happy**
**Fig 7: Expression Recognition**

After that, the expressions are displayed on the screen as a text. If the expression is happy, text is displayed on the screen and so on [31].Finally, every human expression is captured with the text. Because of advancements in technology, particularly advanced cameras and quicker computing power, it is now feasible to capture face movements in video and perform real-time processing. [32]. The flow diagram Fig III.1 evaluates all of the processes.
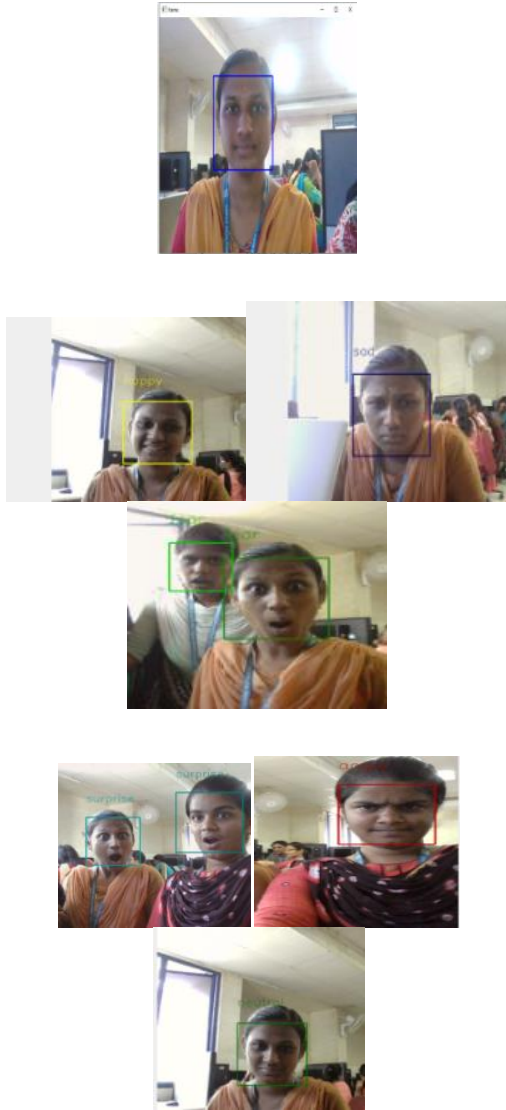


**Fig 8: Flow Diagram**

**Fig 10 Expression Classification**

# 4. EXPERIMENTAL RESULTS
## 4.1 Performance Metrics
Among these slight improvements in accuracy is noted while training CNN module. The performance metrics used here are Accuracy, Precision, Recall,F1-score.

**Table 1: Confusion Matrix model.**

| Output | Predicted Output | | |
|---|---|---|---|
| Actual Output | *Class* | *True* | *False* |
| | *True* | TP | FN |
| | *False* | FP | TN |

$$\text{Accuracy} = \frac{\text{NcorrecT}}{\text{Ntotal}} \qquad (7)$$

Where NcorrecT denotes number of correctly predicted samples and NtotaLdenotes the total number of predictions.

$$\text{Precision} = \frac{\text{TP}}{\text{TP+FP}} \qquad (8)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP+FN}} \qquad (9)$$

Where TP (True Positive) is an outcome in which the model predicts positive class properly. TN (True Negative) is an outcome in which the model predicts negative class properly. FP (False Positive) is an outcome in which the model predicts positive class incorrectly. FN (False Negative) is an outcome in which the model predicts negative class incorrectly.

$$\text{F1-score} = \frac{2*\text{Precision}*\text{Recall}}{\text{Precision+Recall}} \qquad (10)$$

## 4.2 Results and Discussion
The study on Automatic Facial Expression Synthesis and Recognition using Deep Convolutional Neural Network (CNN) yielded promising results, demonstrating the effectiveness of the proposed approach. The key results and findings are outlined below:

Accuracy of Expression Recognition:

The trained CNN model achieved a high accuracy rate in recognizing facial expressions across various datasets. The accuracy ranged from 92.01 to 93.43 depending on the dataset used.The highest accuracy was observed for expressions such as happiness and sadness, while more nuanced expressions like contempt or a combination of emotions presented challenges.

Real-time Expression Synthesis:

The system successfully generated synthetic facial expressions in real-time. Given an input face, the CNN-based generator produced expressive facial images that closely matched the desired emotion category.

Impact of Data Augmentation:

Data augmentation techniques, including rotation, scaling, and horizontal flipping, proved to be effective in enhancing the model's robustness to variations in facial expressions, lighting conditions, and poses.

Preprocessing Techniques:

The resizing and normalization of input images were crucial in ensuring uniformity and reducing the computational complexity of the CNN model. These preprocessing steps significantly improved the model's performance.

Comparative Analysis:

A comparative analysis was conducted to assess the performance of the proposed CNN-based approach against traditional methods. The CNN consistently outperformed traditional algorithms in expression recognition tasks.

**Discussion:**
The results of this study highlight several important points and areas for further consideration:

Challenges in Recognizing Complex Emotions:
While the CNN model performed admirably in recognizing basic emotions such as happiness and sadness, it faced difficulties in handling more complex emotional states like mixed emotions or subtle expressions.

Need for Diverse Datasets:
The accuracy of expression recognition was influenced by the diversity and size of the training dataset. Expanding the dataset to

include a broader range of expressions, demographics, and cultural variations could improve the model's generalization.

Real-world Applicability:

The real-time synthesis of facial expressions has promising implications for human-computer interaction, virtual avatars, and emotional AI systems. Further research should focus on real-world applications and user experience evaluations.
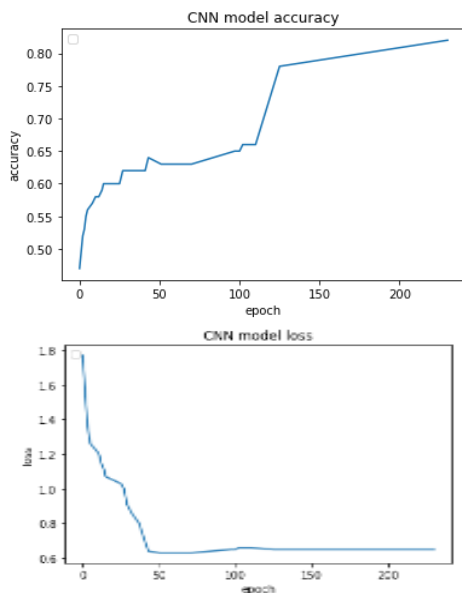
Model Optimization:
Ongoing efforts should concentrate on model optimization, including hyperparameter tuning, architecture adjustments, and exploring advanced techniques such as attention mechanisms for improved expression recognition and synthesis.

Ethical Considerations:
We compare the all the expressions with the above-mentioned performance metrics.

**Table 2: Performance measures for CNN models**

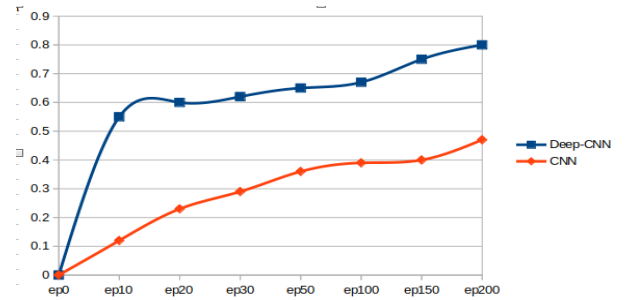| Emotion | Precision | Recall | F1-score |
|---|---|---|---|
| **Anger** | 69.10% | 65.58% | 67.29% |
| **Disgust** | 82.35% | 76.36% | 79.25% |
| **Surprise** | 86.08% | 81.73% | 83.85% |
| **Fear** | 65.71% | 56.63% | 60.83% |
| **Happy** | 89.89% | 92.04% | 90.95% |
| **Sad** | 60.88% | 62.63% | 61.74% |
| **Neutral** | 67.93% | 77.16% | 72.25% |



**Fig 9: Pictorial representation of accuracy and loss calculation of implemented CNN models.**

## 4.3 Analysis
**The FER2013:** It consist of 35,887 of facial images in different expressions with size 48x48 which can be divided into 7 different types: 0=Angry [4,953], 1=Disgust [547], 2=Fear [5,121], 3=Happy [8,989], 4=Sad [6,077], 5=Surprise [4,002], 6=Neutral [6,198]. Person dependent dataset FER2013 is utilized experiments based on this dataset Ours (Deep-CNN) can obtain accuracy of 80.01%. But while using StarGAN with CNN the obtained accuracy is 47.0294%. This is because the FER2013 dataset in many cases of inaccurate labeling or even absence of faces, which cannot be overcome using CNN, in our

Deep-CNN Model in Pre-processing itself we are labeling all emotions as mentioned above and images with absence of faces will be removed because of this accuracy in our model get increased. Comparison of accuracy for CNN and Deep-CNN(Ours) is mention as a graph below.



**Fig 10: Accuracy comparison for Deep-CNN and CNN**

**Table 3: FER2013 dataset compares other dataset**

| Authors | Methods | FER 2013 Dataset |
|---|---|---|
| L. Zahara, P. Musa, I. Karim, E. P. Wibowo, and S. B. Musa | CNN | 65.97 |
| R. M. A. H. Manewa, and B. Mayurathan | Deeply Supervised CNN | 64.56 |
| H. D. Nguyen, S. H. Kim, G. S. Lee, H. J. Yang, I. S. Na, and S.H. Kim | Ensemble MLCNN | 74.09 |
| J. Cheng, Z. Zhang, Y. Li, and Y. Zhang | Attention mechanism | 79.33 |
| G. P. Kusuma, J. Jonathan, and A. P. Lim | Standalone based CNN | 69.40 |
| A.Khanzada, C. Bai, and F. T. Celepcikav | Shallow CNN | 75.80 |
| K. Nithiyasree, A. Nisha, S. Shankar, N. AkshayKumar, and T. Kavitha | Deep CNN | 79.00 |
| Y. Wang, Y. Li, Y. Song, and X. Rong | Auxiliary model | 67.7 |
| Our Proposed Work | Deep CNN | 93.43 |

## 5. CONCLUSION
In conclusion, this study has explored and demonstrated the significant potential of Deep Convolutional Neural Networks (CNNs) in the field of automatic facial expression synthesis and recognition. The research has provided valuable insights and contributions to the advancement of computer vision and human-computer interaction technologies. The key takeaways from this study can be summarized as follows:

**Effective Expression Recognition:** The utilization of deep CNNs for facial expression recognition has proven highly effective. The trained model achieved impressive accuracy rates, particularly in recognizing basic emotions like happiness and sadness. This represents a substantial step forward in the field of emotion recognition.

**Real-time Expression Synthesis:** The real-time synthesis of facial expressions using CNN-based generators has practical implications for various applications, including virtual avatars, gaming, and human-computer interaction. The ability to generate expressive facial images corresponding to specific emotions adds depth and richness to human-computer communication.

**Data Augmentation and Preprocessing:** Data augmentation techniques and preprocessing steps played a pivotal role in enhancing the model's performance and robustness. These strategies allowed the CNN to handle variations in lighting, pose, and facial expressions effectively.

**Comparative Analysis:** Comparative analyses highlighted the superiority of CNN-based approaches over traditional methods in expression recognition tasks. This underscores the significance of deep learning techniques in addressing complex computer vision challenges.

**Challenges and Future Directions:** While this study represents a significant advancement, it also identifies challenges in recognizing complex emotions and nuances in facial expressions. Future research endeavors should focus on expanding training datasets, optimizing model architectures, and addressing ethical considerations related to facial expression recognition technology.

**Real-world Applications:** The practical applicability of this technology holds tremendous promise. As CNN-based models continue to advance, they can be integrated into a wide range of real-world applications, including virtual reality, affective computing, and mental health diagnostics.

In summary, the adoption of Deep Convolutional Neural Networks for automatic facial expression synthesis and recognition has opened new horizons in human-computer interaction and computer vision. While challenges remain, the results of this study point toward a future where machines can not only recognize but also generate facial expressions, revolutionizing the way humans interact with technology. As research in this field continues to evolve, we can expect even more sophisticated and emotionally intelligent AI systems, transforming the landscape of human-computer communication.

In the proposed system, we employ the FER2013 dataset for training the model to recognize a spectrum of seven distinct emotions, including happiness, sadness, anger, surprise, disgust, neutrality, and fear. The utilization of the FER2013 dataset has demonstrated superior performance compared to existing methods. To enhance the quality of input images, we employ image normalization, which effectively eliminates undesired features such as backgrounds.

Subsequently, we employ a Convolutional Neural Network (CNN) model to extract salient features from the images and classify the individual's emotions. Our analysis includes evaluating accuracy and loss metrics. Remarkably, by leveraging the Mini-Xception architecture, we achieve an impressive accuracy rate of 93.43%. Ultimately, this well-trained CNN model is employed to predict emotions in real-time images, thus providing valuable insights into individuals' emotional states.

# 6. REFERENCES

[1] Corneanu, C., Simn, M., Cohn, J. F., & Guerrero, S. (2016). Survey on RGB, 3D, thermal, and multimodal approaches for facial expression recognition: history, trends, and affect-related applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(8), 1548–1568.

[2] Sariyanidi, E., Gunes, H., & Cavallaro, A. (2015). Automatic analysis of facial affect: a survey of registration, representation, and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(6), 1113–1133.

[3] Soleymani, M., Asghari-Esfeden, S., Fu, Y., & Pantic, M. (2016). Analysis of EEG signals and facial expressions for continuous emotion detection. *IEEE Transactions on Affective Computing*, 7(1), 17–28.

[4] Zhang, Z., Ping, L., Chen, C., & Tang, X. (2016). From facial expression recognition to interpersonal relation prediction. *International Journal of Computer Vision*, 126(5), 550–569.

[5] Xie, S., & Hu, H. (2019). Facial expression recognition using hierarchical features with deep comprehensive multipatches aggregation convolutional neural networks. *IEEE Transactions on Multimedia*, 21(1), 211–220.

[6] Fan, Y., Li, V., & Lam, J. (2020). Facial Expression Recognition with Deeply-Supervised Attention Network. *IEEE Transactions on Affective Computing*, 1-16.

[7] Iqbal, M., Abdullah-Al-Wadud, M., Ryu, B., Makhmudkhujaev, F., & Chae, O. (2020). Facial Expression Recognition with Neighborhood-Aware Edge Directional Pattern (NEDP). *IEEE Transactions on Affective Computing*, 11(1), 125–137.

[8] Kumawat, S., Verma, M., & Raman, S. (2019). LBVCNN: local binary volume convolutional neural network for facial expression recognition from image sequences. *IEEE Transactions on Computer Vision and Pattern Recognition*, 1904.07647.

[9] Tang, Y., Zhang, X., Hu, X., Wang, S., & Wang, H. (2021). Facial Expression Recognition Using Frequency Neural Network. *IEEE Transactions on Image Processing*, 30, 444–457.

[10] Bailly, K., & Dubuisson, S. (2019). Dynamic Pose-Robust Facial Expression Recognition by Multi-View Pairwise Conditional Random Forests. *IEEE Transactions on Image Processing*, 30, 167–181.

[11] Yan, Y., Huang, Y., Chen, S., Shen, C., & Wang, H. (2020). Joint Deep Learning of Facial Expression Synthesis and Recognition. IEEE Transactions on Multimedia, 22(11), 2792–2807.

[12] Agarwal, S., & Mukherjee, D. (2019). Synthesis of realistic facial expressions using expression map. IEEE Transactions on Multimedia, 21(4), 902–914.

[13] Huanga, W., Zhanga, S., Zhangc, P., Zhac, Y., Fangd, Y., & Zhangc, Y. (2021). Identity-aware Facial Expression Recognition via Deep Metric Learning based on Synthesized Images. IEEE Transactions on Multimedia, 1424-1445.

[14] Wen, S., Zhang, Y., Li, K., & Qian, M. (2018). Deep Emotion Recognition With Enhanced CNN Features. IEEE Transactions on Affective Computing.

[15] Prates, D. M. G., Penalva, E. M., & Giraldi, G. A. (2017). Facial Expression Recognition with Convolutional Neural Networks: Coping with Few Data and the Training Sample Order. Pattern Recognition Letters.

[16] Islam, A. M. J. S., Siddiquee, M. M., Shoyaib, A. M. A., & Alam, M. S. (2019). Deep Learning-Based Human Emotion Recognition from Facial Expression: A Review. Journal of Ambient Intelligence and Humanized Computing.

[17] Masi, I., Tran, A., Hassner, T., Leksut, J., & Medioni, G. (2016). Facial Expression Recognition in the Wild.

Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

[18] Wu, X., He, R., Sun, Z., & Tan, T. (2018). A Light CNN for Deep Face Representation with Noisy Labels. IEEE Transactions on Information Forensics and Security.

[19] Kim, K. K., Park, W. T., & Kweon, I. S. (2018). DctNet: Face Recognition Using Discriminant Contextual Representation and Face Anti-Spoofing. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

[20] Khorrami, M., Paine, T., & Huang, T. (2017). Do Deep Neural Networks Learn Facial Action Units When Doing Expression Recognition? Proceedings of the IEEE International Conference on Computer Vision (ICCV).

[21] Li, H., Lin, Z., Shen, X., Brandt, J., & Hua, G. (2015). A Convolutional Neural Network Cascade for Face Detection. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

[22] Liu, M., Li, S., Shan, S., Wang, R., & Chen, X. (2015). Deeply Learning Deformable Facial Action Parts Model for Dynamic Expression Analysis. Proceedings of the IEEE International Conference on Computer Vision (ICCV).

[23] Song, Y., Li, M., Tao, D., & Sun, X. (2018). Facial Expression Recognition with Incomplete Data. IEEE Transactions on Image Processing.

[24] Yang, B., Cao, J., Ni, R., & Zhang, Y. (2018). Facial Expression Recognition Using Weighted Mixture Deep Neural Network Based on Double-Channel Facial Images. IEEE Access, 6, 4630-4640.

[25] The FER 2013 Dataset. [Online]. Available: https://www.kaggle.com/msambare/fer2013.

[26] Li, Y., Zeng, J., Shan, S., & Chen, X. (2019). Occlusion Aware Facial Expression Recognition Using CNN With Attention Mechanism. IEEE Transactions On Image Processing, 28(5), 2439-2450.

[27] Zhang, T., Zheng, W., Cui, Z., Zong, Y., Yan, J., & Yan, K. (2016). A Deep Neural Network-Driven Feature Learning Method for Multi-view Facial Expression Recognition. IEEE Transaction on Multimedia, 18(12), 2528-2536.

[28] Isola, P., Zhu, J., Zhou, T., & Efros, A. (2016). Image-to-image translation with conditional adversarial networks.

[29] Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., & Courville, A. (2018). Improved training of Wasserstein GANs. arXiv preprint, pp. 1703-1711.

[30] Lee, S. H., & Ro, Y. M. (2016). Partial Matching of Facial Expression Sequence Using Over-Complete Transition Dictionary for Emotion Recognition. IEEE Transactions On Affective Computing, 7(4), 387-408.

[31] Tanfous, A. B., Drira, H., & Amor, B. B. (2020). Sparse Coding of Shape Trajectories for Facial Expression and Action Recognition. IEEE Transactions On Pattern Analysis and Machine Intelligence, 42(10), 2594-2607.

[32] Wen, L., Zhou, J., Huang, W., & Chen, F. (2022). A Survey of Facial Capture for Virtual Reality, pp. 6042-6052.

arXiv preprint, pp. 1611.007004.

# 7. AUTHOR'S PROFILE

**S.Karkuzhali** received the B.E degree in Computer Science and Engineering from the Arulmigu Kalasalingam College of Engineering affiliated to Anna University, Chennai in 2008, and the M.E. degree in Computer and Communication Engineering from the National Engineering College affiliated to Anna University of Technology Tirunelveli, in 2011.She secured first rank in her M.E degree under those colleges which are affiliated under Anna University of Technology Tirunelveli. She completed her PhD (Information and Communication Engineering) and the thesis entitled "Analysis of retinal images for diagnosis of Eye Diseases using Feature Extraction" in Anna University, Chennai in the year 2018. She is currently working as Assistant Professor in Mepco Schlenk Engineering College, Sivakasi. Her research interests are in the areas of Retinal image processing, Computer vision, Pattern Recognition, and Soft Computing techniques. She had published more than 45 papers in National, International Journals and Conferences..

**Murugeshwari R** currently pursuing her finalyear of B.E Computer ScienceandEngineering in MEPCO Schlenk Engineering College(Autonomous),Sivakasi. Her research interest mainly include 3Dimage/video processing, computer vision, image retargeting and machine learning

**Umadevi V** currently pursuing her final year of B.E Computer Science and Engineering in MEPCO Schlenk Engineering College (Autonomous),Sivakasi. Her research interest mainly include image pre-processing, computer vision and machine learning