

Original Article

# Emotion Recognition Using PIZAM-ANFIS by Considering Partial Occlusion and Behind the Mask

Jyoti S. Bedre<sup>1\*</sup>, P. Lakshmi Prasanna<sup>1</sup>

<sup>1</sup>Computer Science and Engineering, KL University, Andhra Pradesh, India.

\*Corresponding Author : [jyoti.phd2020@gmail.com](mailto:jyoti.phd2020@gmail.com)

Received: 08 June 2024

Revised: 10 January 2025

Accepted: 07 February 2025

Published: 21 February 2025

**Abstract** - Emotional expressions, encompassing verbal and non-verbal communication, convey an individual's emotional state or attitude to others. Understanding complex human behavior requires analyzing physical features across multiple modalities, with recent studies focusing extensively on spontaneous multi-modal emotion recognition for human behavior analysis. However, accurate Facial Emotion Recognition (FER) faces significant challenges due to partial facial occlusions caused by random objects and mask-wearing. The paper introduces a novel classification method, Pizam-ANFIS-based FER, which considers Occlusions and Masks (PAFEROM) to address this. Preprocessing the input image is the first step in the process, followed by cropping and face detection with the Viola-Jones Algorithm (VJA). Further, the skin tone is then analyzed, and several parts of the face are segmented using LSW-KCM. Furthermore, contour formation, edge detection by CGED, and extracting features are executed. Using Principal Component Analysis with Information Gain Analysis (PIGA), the retrieved features dimensionality is reduced before the CSE processes them for the identification of Action Units (AUs), and the proposed approach is utilized. Subsequently, the identified AUs and dimensionally reduced features are classified using Pizam-ANFIS to recognize human emotions. Experimental results indicate that the proposed model surpasses existing techniques in both effectiveness and accuracy.

**Keywords** - Local Structural Weighted K-Means Clustering (LSW-KMC) algorithm, Canny Gaussian Edge Detector (CGED), PizMamdani (Pizam)-Adaptive Neuro Fuzzy InterferenceSystem (Pizam-ANFIS), Correlated Swish Embedding Network (CSE).

## 1. Introduction

Facial expressions play an important role in human communication by providing essential nonverbal information that complements verbal interactions. Studies suggest that a significant portion of communication, ranging from 60% to 80%, is conveyed through nonverbal cues. These include facial expressions, eye contact, vocal tone, hand gestures, and physical distance [1, 2]. Analyzing these facial expressions has garnered significant attention in research, particularly in the field of FER. FER technology is increasingly utilized in Human-Computer Interaction (HCI) applications, including autopilot systems, education, medical and psychological treatments, surveillance, and psychological analysis in computer vision [3]. By examining human facial expressions, FER systems aim to detect specific emotions such as anger, disgust, fear, happiness, sadness, surprise, and neutral states. The complexity of accurately estimating emotions is heightened by the diversity of human facial features and the variety of possible emotional expressions [4]. Automated recognition of facial expressions has garnered significant interest in recent years due to its broad spectrum of applications [5]. However, achieving high accuracy in recognizing facial expressions remains challenging because of

their subtlety, complexity, and diversity [6]. A critical aspect of effective FER is obtaining precise facial representations from the original facial images [7].

This system has two tasks: face detection and facial emotion classification. To extract significant and unique facial features, the human face is first recognized from the acquired image [8]. Then, the emotion represented by the identified face is classified using a FER algorithm. Formerly, researchers have tackled FER using various approaches such as the MLP Model, k-nearest Neighbors (KNN), and Support Vector Machines (SVM) [9] have been used to extract information through methods such as Local Binary Patterns, Eigenfaces, Face-Landmark, and Texture features. Among these approaches, neural networks have gained significant popularity and are now widely utilized for FER [10]. Presently, advanced classifiers, including Artificial Neural Networks (ANN), Convolutional Neural Networks (CNN), and Random Forests, are extensively employed in this domain and are widely used for tasks such as healthcare recognition, biometric identification, handwriting analysis, and facial detection for security purposes. However, achieving precise emotion classification with state-of-the-art classifiers in FER



remains challenging due to issues like partial occlusion and the use of masks, which often need to be adequately addressed.

### 1.1. Problem Statement

Listed below are some of the shortcomings of the existing research approaches used to date:

1. Although current facial expression classifiers have proven practically flawless in analyzing confined frontal faces, there is a need to improve when analyzing faces that are partially obscured or hidden behind masks, frequently seen in the wild.
2. When wearing a face mask that covers the mouth and nose, it is impossible to accurately identify facial expressions of emotion. Classifying facial emotions using the half-face is more complex and challenging since the mouth area is one of the significant variables responsible for emotion detection.
3. Current FER techniques for masked faces often disregard significant facial areas like the forehead. Instead, they isolate only the eye region using landmark detection methods, which ultimately reduces the accuracy of the FER system.

The limitations of traditional FER systems become particularly evident in real-world scenarios involving partial occlusions and mask-wearing, which have become increasingly common in recent years. For instance, in healthcare settings, masked faces of medical personnel pose challenges for FER systems attempting to monitor stress or fatigue levels. Similarly, in surveillance applications, occluded faces due to scarves, helmets, or other coverings often result in inaccurate emotion detection, potentially undermining security measures. In education, where FER is used to assess student engagement during online learning, using face masks or partial visibility due to camera angles can lead to misclassification of emotions, reducing the efficacy of such tools. These real-world challenges underscore the need for robust FER systems capable of accurately detecting emotions under such conditions. This research addresses these limitations by introducing a novel framework designed to maintain high accuracy even in occlusion-prone scenarios, as demonstrated by its superior performance compared to existing techniques.

This research suggests an improved FER system using a novel Pizam-ANFIS classifier to overcome these issues. The key research objectives of this system are outlined as follows:

1. A novel Edge Detector has been developed to detect the exact boundaries of organs.
2. A novel dimensionality reduction model is employed to select the interest features to mitigate training time.
3. A novel neural network is employed
4. To categorize the AU present in the mask-covered facial image.

5. A rule-based novel technique is utilized to classify human emotions.

Despite significant advancements in FER, existing techniques face notable challenges in accurately classifying emotions in real-world scenarios involving partial occlusion, such as mask-wearing or object obstruction. Addressing these gaps, this study introduces the Pizam-ANFIS classifier. This novel framework integrates advanced edge detection, dimensionality reduction, and action unit identification techniques to enhance FER performance under challenging conditions. Unlike prior approaches, the proposed model utilizes LSW-KMC for precise feature extraction and a PIGA-based dimensionality reduction method to optimize computational efficiency.

The novelty of this work lies in its ability to achieve superior accuracy and robustness, particularly in scenarios involving occlusions, which outperform existing models. The outline of this paper is as follows: Section 2 offers an in-depth review of related work, emphasizing significant advancements and challenges within the field. Section 3 details the proposed methodology, highlighting the innovative techniques and algorithms employed. Section 4 then presents and analyzes the results, emphasizing performance metrics and comparative evaluation. At last, the 5<sup>th</sup> Section concludes the paper by summarizing the findings and suggesting potential directions for future research.

## 2. Literature Survey

In the field of FER, Mehendale et al. [11] introduced a modular framework that employs an AdaBoost cascade classifier for face detection and utilizes Neighborhood Difference Features (NDF) for feature extraction, which were then classified using a random forest classifier to address false detections. Despite outperforming methods on the SFEW and RAF datasets, the system's omission of geometric elements led to inaccuracies. Liu et al. [12] introduced a FER technique that utilized landmark curvature and vectorized landmarks, blending SVM classification with a GA to select features and parameters.

While this approach showed balanced performance on the CK+ and MUG datasets, image noise impacted the SVM classifier's accuracy. Alreshidi et al. [13] employed NPCA for dimensionality reduction and SVM for emotion recognition, achieving high accuracy but struggling with varying input dimensions. Hassan et al. [14] utilized graph mining techniques to identify common sub-graphs within emotional classes, enhancing efficiency and accuracy but resulting in a more time-consuming process. Hussain et al. [15] developed a deep learning-based FER system structured in three phases: face detection, feature analysis using Keras CNN, and emotion classification. Although this system demonstrated proficiency, errors in facial landmark detection impacted overall accuracy.

Houshmand et al. [16] proposed a transfer learning approach with pre-trained VGG and ResNet networks for FER under VR headset occlusion, achieving comparable performance but needing refinement in preprocessing steps due to issues with histogram equalization. Monisha et al. [17] introduced a real-time FER system using CNN for classification, demonstrating high accuracy but encountering recognition errors due to limited training data. Akhand et al. [18] utilized transfer learning within a Deep Convolutional Neural Network (DCNN), progressively enhancing FER accuracy but failing to preserve edge information crucial for detailed emotion recognition. Saha et al. [19] employed the Cosine Similarity-Based Harmony Search Algorithm (SFHSA) for feature selection, optimizing feature vectors and improving classification accuracy, albeit with a time-consuming training process. Gautam et al. [20] combined HOG and SIFT for extracting features with classification using CNN, outperforming existing methods but struggling with the limitations of 2D data in handling facial pose variations. Castellano et al. [21] focused on recognizing emotions from masked faces using ResNet, achieving high accuracy with eye region analysis but increasing computational demands due to skipping connections.

Wally et al. [22] developed an Occlusion-Aware Student Emotion Recognition system utilizing CNN and FCNN, which faced overfitting issues due to limited data. Elsayed et al. [23] showcased a hybrid CNN with LBP for feature extraction in masked faces, demonstrating improved recognition but facing challenges with imbalanced and noisy data. Mukhiddinov et al. [24] applied synthetic masks to input images, emphasizing head and upper facial features for FER, achieving higher accuracy but encountering orientation issues with landmark features. Finally, Zhu et al. [25] introduced HDCNet, leveraging a feature constraint methodology to mine attention consistency features, improving classification accuracy but posing substantial computational demands due to Class Activation Mapping.

### 3. Proposed Framework for FER

This study introduces a novel Pizam-ANFIS model for accurate and efficient human emotion recognition using visual features. Two key processes face detection and classification are finished in order to identify the facial mood. Features from the face are retrieved and fed into a trained network for emotion classification. The block diagram for the suggested model is illustrated in Figure 1.

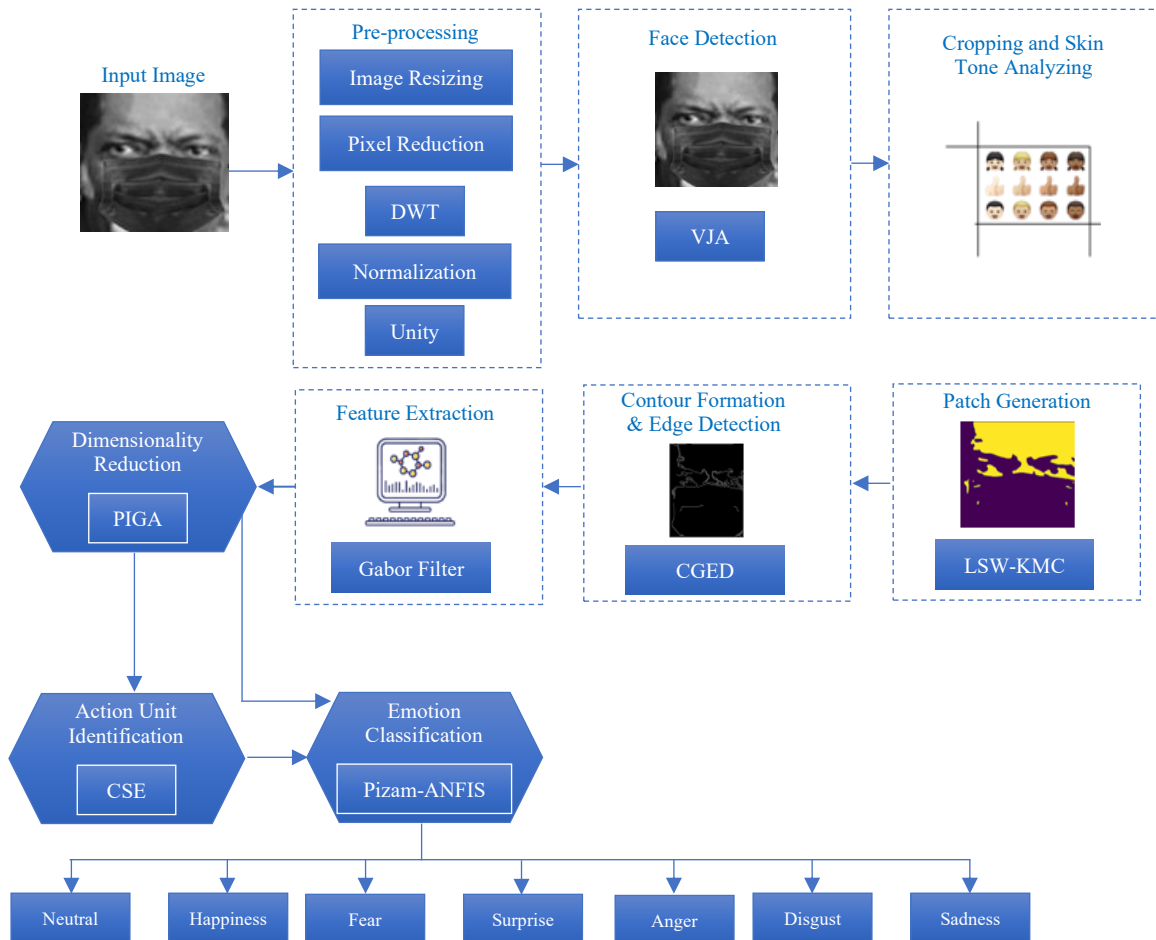


Fig. 1 Schematic of the projected framework

### 3.1. Preprocessing

This section demonstrates how an image is initially taken as input and processed through preprocessing to eliminate unwanted elements. The input face expression image undergoes the preprocessing operation in three stages: Image resizing, pixel reduction and normalization.

#### 3.1.1. Image Resizing

The accuracy and computation time of the processing system can be adversely affected by unwanted pixels in the input image. The input image (I) is resized to 256×256 pixels using bilinear interpolation to address this. This method is particularly recommended for continuous data sets lacking distinct boundaries. Bilinear interpolation is a resampling method that computes a new pixel value by averaging the four nearest pixel values, weighted by their distances. This technique provides a smoother and more precise representation of the image. The resized image ( $I_{resize}$ ) is,

$$I_{resize} = \frac{\psi^R S^L + \psi^L S^R + \psi^T S^B + \psi^B S^T}{\psi^R + \psi^L + \psi^T + \psi^B} \quad (1)$$

Here,  $\psi^R, \psi^L, \psi^T,$  and  $\psi^B$  refer to the distances corresponding to the missing pixel. And  $S^L, S^R, S^B,$  and  $S^T$  represents the source pixels located to the left, right, top, and bottom.

#### 3.1.2. Pixel Reduction

After resizing the image, the noisy pixels from the resized image ( $I_{resize}$ ) were removed by utilizing the Discrete Wavelet Transform (DWT). DWT is selected due to its ability to achieve a higher compression ratio. This process involves decomposing the image into coefficients (sub-bands) and then compared to a set threshold ( $T_{thres}$ ). The coefficients that fall below this threshold are considered noiseless pixels and are retained in the image. In contrast, those above the threshold are identified as noisy pixels and are subsequently removed. This method ensures that only the low-low frequency sub-bands, which contain the essential image information with reduced noise, are preserved. The resulting pixel-reduced image ( $I_{red}$ ) can be represented as follows:

$$I_{red} = \begin{cases} \text{noisy pixel, if } (I^{red})^\rho > T_{thres} \\ \text{noiseless pixel, if } (I^{red})^\rho < T_{thres} \end{cases} \quad (2)$$

#### 3.1.3. Normalization

Unity normalization transforms the pixel-reduced image ( $I^{red}$ ) into a range of pixel values. Unity normalization has better and faster execution. In order to reduce the inner-class feature mismatch, which can be seen as intensity offsets, image normalization is a crucial preprocessing approach. The normalized image can be denoted as  $I^{nor}$ .

$$I^{nor} = \frac{I^{red}}{\|v\|} \quad (3)$$

Here,  $\|v\|$  denotes the vector of the pixels.

### 3.2. Face Detection

In this step, face detection from the preprocessed image using the VJS to facilitate the determination of the region of interest and subsequent feature extraction. The VJS process entails sliding feature boxes across the image and computing the difference in the total pixel values between adjacent regions, represented as (d). This difference is then compared to a threshold value ( $T_f$ ) to determine if an object, such as a face, has been detected. This method simplifies the identification of the region of interest and ensures accurate feature extraction from the detected face. The detected face ( $I^{face}$ ) is computed as follows,

$$T_f = \begin{cases} \text{face detected if } T_f > d \\ \text{not detected if } T_f < d \end{cases} \quad (4)$$

### 3.3. Cropping and Skin Tone Analysis

The detected face image is cropped to remove all the unwanted things from the image, such as the background, and to keep only relevant information in the image. After that, skin tone analysis is done to differentiate the parts presented over the face. Then, the image of the skin analyzed is denoted as  $I^{skin}$ .

### 3.4. Patch Generation

In patch generation, the different facial parts are segmented from  $I^{skin}$  to encourage extracting discriminative features from the minute parts using the LSW-KMC algorithm. The LSW-KMC algorithm was employed for precise segmentation of facial regions. This method utilizes a weighted sum of image pixels to improve clustering accuracy, with the control parameter  $\alpha$  set to balance spatial and structural relationships. The structural similarity index, combining luminance, contrast, and structural metrics, guides the clustering process until convergence. K means is favored over other segmentation methods because of its ease of use and rapid computation speed. However, the spatial Euclidean distance-based characterization of the relationship between the image pixels and cluster center is more difficult since this distance alone is insufficient to understand the general characteristics. In order to get over the drawbacks above, the weighted sum of the image pixels was used to estimate the distance between each image pixel and the cluster center. After that, the structural similarity index calculates a local distance measurement to determine how far apart two image pixels are from one another in the overall image. This local distance computation reflects not only the physical relationship between two picture pixels but also the relationship connected to luminance and contrast, as well as the structure of the image pixels revolving around them. As a result, LSW-KMC serves as the inspiration for the proposed KMC. The steps of LSW-KMC are listed as:

a) Initializing the pixels  $\rho^j \in I^{skin}$ , presented as,

$$\rho^j = \{\rho^1, \rho^2, \rho^3, \dots, \rho^N\} \text{ where, } j = 1, 2, 3, \dots, N \quad (5)$$

Here,  $j$  denotes the count of pixels of the skin tone detected image.

- b) Select the cluster numbers that are defined by respective centroids. Initially, the precise centers of the pixels are unknown, so the centroids  $C_m$  are chosen randomly to establish each cluster.

$$C_i = C_1, C_2, \dots, C_M \quad i = 1, 2, \dots, M \quad (6)$$

Here,  $i$  represents the centroid (cluster centre).

- c) Calculate the weighted sum of the image pixels by considering the essential distance  $(d(\rho^j, C_i))$ .

$$S = \sum_{r=1}^Z W_R d(\rho^{jr}, C_{ir}) \quad (7)$$

Here,  $W_R$  denotes the weight associated with the distance  $(d(\rho^j, C_i))$ ,  $\rho^{jr}$  represents the value of the point in the image located around the  $\rho^j$ ,  $C_{ir}$  denotes centroids, and  $Z$  denotes the number of points in the skin tone detected image.

- d)  $W_R$  is determined by looking at the coordinate distance between  $\rho^{jr}$  and  $\rho^j$ . Therefore, the weights are,

$$W_R = \frac{1}{(1+d_r)^{C_{para}}} \quad (8)$$

Here,  $C_{para}$  represents the control parameter.

- e) Measure the structural similarity of the image. It considers the degree of similarity of luminance, contrast, and structure of the pixel and cluster center. The SSIM index  $(D \in S)$  between pixels and cluster center is defined as,

$$D = \frac{(2\lambda_{\rho^j} \lambda_{C_i} + \chi_1)(2\sigma_{\rho^j C_i} + \chi_2)}{(\lambda_{\rho^j}^2 \lambda_{C_i}^2 + \chi_1)(\sigma_{\rho^j}^2 \sigma_{C_i}^2 + \chi_2)} \quad (9)$$

Here,  $\lambda_{\rho^j}$  and  $\lambda_{C_i}$  denote the mean of  $\rho^j$  and  $C_i$  respectively,  $\sigma_{\rho^j C_i}$  signifies the cross-correlation between  $\rho^j$  and  $C_i$ ,  $\sigma_{\rho^j}^2$  and  $\sigma_{C_i}^2$  specifies the standard deviation of  $\rho^j$  and  $C_i$ , respectively,  $\chi_1$  and  $\chi_2$  are the positive constants.

- f) Assign each pixel to the cluster whose centroid is closest, minimizing the distance between the pixel and the centroid.

This process continues iteratively until the clusters stabilize and no further changes occur. This segmentation identifies and outlines standard and disease-affected regions in the resulting image, denoted as  $I_{seg}$ . The pseudocode for the proposed LSW-KMC means is:

```

Input  : Face-detected image  $I^{skin}$ 
Output : Segmented image  $I^{seg}$ 
Begin
    Initialize  $\rho^n$ , number of clusters  $C_m$ , iteration ( $iter$ ), maximum iteration ( $iter_{max}()$ )
    Perform clustering
    Select the number of centroids
    Set  $iter = 1$ 
    While  $iter \leq iter_{max}$ 
        For each pixel, do
            Calculate the weighted sum of image pixels
            Compute distance  $D$ 

$$D = \frac{(2\lambda_{\rho^i} \lambda_{C_i} + \chi_1)(2\sigma_{\rho^j C_i} + \chi_2)}{(\lambda_{\rho^i}^2 \lambda_{C_i}^2 + \chi_1)(\sigma_{\rho^i}^2 \sigma_{C_i}^2 + \chi_2)}$$

            End for
        Check all the pixels are presented under the cluster
        If ( $\rho^n == undercluster$ ) {
            Stop criteria
        }
        Else
        {
            Set  $iter = iter + 1$ 
        }
        End if
    End While
    Return segmented image
End
    
```

### 3.5. Contour Formation and Edge Detection

Here, the contour is formed over  $I^{seg}$  using CGED to extract the facial parts more effectively from the occluded and mask-covered input images. For simplicity, the existing Canny Edge Detection (CED) is chosen for the proposed work. However, a drawback of the CED is that the default Sobel Operators are restricted to a fixed 3-by-3 window. This limitation can be problematic, particularly in noisy images, potentially compromising the final output. The work employs a broader 5-by-5 Sobel Operator window to address this issue. Additionally, the horizontal and vertical gradients are calculated using a Gaussian kernel, replacing the standard convolution kernel used in traditional CED. This adjustment reduces computational time while enhancing noise resistance and edge detection accuracy, making the CGED approach more robust and effective for occluded and mask-covered facial images. Denoise image before detecting the edge of the image usually use the 5-by-5 Sobel Operator to reduce noise, according to (10),

$$I^{den} = \sqrt{\beta_o^2 + \beta_t^2} \quad (10)$$

To calculate the gradient intensity (B), use the Gaussian kernel and determine the edge direction ( $\phi$ ). Typically, the gradient direction is categorized into four angles: 0, 45, 90,

and 135 degrees. This process is defined by Equations (11) and (12),

$$B = \exp\left(\frac{-\|\beta_o - \beta_l\|^2}{2\sigma^2}\right) \quad (11)$$

$$\phi = \tan^{-1}\left(\frac{\beta_o}{\beta_l}\right) \quad (12)$$

Where,  $\beta_o$  and  $\beta_l$  denote the pixel values in the  $o$ -axis and  $l$ -axis, respectively,  $\sigma$  denotes the signum function. After the gradient and magnitude calculation, the entire image is scanned, unwanted pixel intensities are suppressed to 0, and the edges present are given as  $E_h, h = 1, 2, \dots, fin$ . Next, the hysteresis threshold is selected as high ( $Up_l$ ) and low ( $Lo_l$ ). These thresholds analyze whether all the detected edges are edges or not. The thresholding function is given as,

$$I^{edge} = \begin{cases} \text{Sure edge if } h > Up_l \\ \text{Valid edge if } Up > h > Lo_l \\ \text{non edge else} \end{cases} \quad (13)$$

Where  $h$  depicts the edge, if the edge  $h$  lies between, then  $Up_l$  and  $Lo_l$  connected to a sure edge is considered a valid edge. If the edge  $h$  does not connect to the sure edges and below, then  $Lo_l$  it is removed from the image as a non-edge. Finally, the edge-detected image is denoted as  $I^{edge}$ .

### 3.6. Feature Extraction

After performing edge detection, the next step is to extract features to obtain detailed information from the input image. Texture features are extracted using the GF, a linear filter selected for its frequency and orientation representations that closely mimic the human visual system. The Gabor Filter (GF) comprises a sinusoidal plane wave modulated by a Gaussian kernel function. Based on the convolution theorem, the Fourier Transform (FT) of a harmonic function and the FT of a Gaussian function combine to produce the impulse response of a Gabor filter. This filter captures orthogonal directions with both real and imaginary components. The process involves applying the GF to the input image to obtain the sinusoidal plane wave response, modulating this response with the Gaussian kernel function to capture both frequency and orientation information, and combining the Fourier transforms of the harmonic and Gaussian functions to generate the GF's impulse response. The real and imaginary components representing orthogonal directions are then extracted. These Gabor features ( $f_1$ ) are crucial for accurately capturing the texture information from the image, thereby enhancing the overall feature extraction process.

$$f_1 = \exp(-(\rho^i)^2 + (\rho^i)^2/2\varpi^2) * \cos(2\pi/\lambda) \rho^i \quad (14)$$

Here,  $\lambda$  and  $\varpi$  denotes the wavelength and effective width, respectively. Additionally, various features such as geometrical features, appearance features, temporal features,

HOG, SIFT, and Speeded-Up Robust Features (SURF) are extracted. The comprehensive set of extracted features ( $f_k$ ) can be summarized as follows:

$$f_k = \{f_1, f_2, f_3, \dots, f_K\} \text{ where, } k = 1, 2, 3, \dots, K \quad (15)$$

Here,  $K$  denotes the number of features.

### 3.7. Dimensionality Reduction

In this step, the dimensionality of features is reduced  $f_k$  to a lower-dimensional space using PIGA, which selects the most critical features to minimize training time during classification. Principal Component Analysis (PCA) is employed for its straightforward computation process and ability to eliminate correlated features. Principal Components aim to capture the maximum variance among the features. However, traditional PCA may lose some information compared to the original feature set due to the arbitrary selection of principal components. To address this limitation, the research incorporates the Information Gain (IG) mechanism, an entropy-based feature estimation method, to determine the optimal number of principal components. IG evaluates each feature individually, calculates its information gain, and assesses its importance concerning the class label.

Each extracted feature is assigned a score ranging from 1 to 0, indicating its relevance from most to least important for setting the number of principal components. The covariance matrix for PCA is computed using the normalized features, and eigenvalues are calculated using decomposition functions, which are then ranked based on IG scores to prioritize the most significant features. For this study, the threshold for IG was experimentally set at 0.5 to ensure a balance between feature retention and dimensionality reduction. By combining PCA's ability to optimize variance with IG's feature evaluation, PIGA ensures that principal component selection is fair and effective and preserves essential information while reducing dimensionality.

#### 3.7.1. Covariance Matrix Construction

The PIGA constructs a covariance matrix for the recognition process to get the eigenvectors. The covariance matrix ( $\mathfrak{R}$ ) construction is formulated as,

$$\mathfrak{R} = \frac{1}{K} \sum_{k=1}^K (f_k) (f_k)^T \quad (16)$$

Where,  $(f_k)^T$  depicts matrix transpose.

#### 3.7.2. Eigenvalue Calculation

The eigenvalue is calculated from the features as,

$$E = \vartheta((1/K) \times f_k) \quad (17)$$

Where  $E$  depicts the eigenvalue and  $\vartheta((f_k)^T)$  depicts the decomposition function, which is given as,

$$\vartheta = D_{co}D_{main} \quad (18)$$

Here,  $D_{co}D_{main}$  depicts the decomposition of two matrices of the features.

### 3.7.3. Eigenvector Estimation

For the features with high eigenvalues, the eigenvector (V) is calculated using the formula,

$$V = \mathfrak{R} - \zeta . E \quad (19)$$

Here,  $\zeta$  indicates a random constant value.

### 3.7.4. Obtaining Principal Components

After the eigenvalues are estimated, the features with high Eigenvalues are derived as the principal components. The principal components are calculated using IG,

$$p_{com} = V \times \varphi_{cen} \quad (20)$$

Where,  $\varphi_{cen}$  depicts the kernel center. Thus, the selected features ( $F_z^{sel}$ ) are given as,

$$F_z^{sel} = [F_1^{sel}, F_2^{sel}, F_3^{sel}, \dots, F_z^{sel}] \quad (21)$$

Where  $Z^{th}$  represents the number of features.

The methodology is designed as a cohesive pipeline where each stage builds on the outputs of the previous one to achieve accurate emotion recognition. Using the LSW-KCM algorithm, the segmentation stage extracts specific facial regions by clustering pixels based on spatial and structural similarities. This segmentation ensures that key facial features are isolated for further processing. The segmentation output is passed to the edge detection stage, where the CGED method is employed. By using a broader Sobel operator and Gaussian kernel, CGED effectively refines the contours of segmented regions, ensuring precise boundary detection even in noisy or occluded images. These refined edges provide a robust input for the feature extraction stage, where critical texture and structural features are identified using techniques like Gabor filters. The extracted features are then fed into the dimensionality reduction stage. PIGA selects the most informative features based on variance and information gain, reducing computational complexity while retaining essential data. This sequential integration ensures that each stage enhances the quality and relevance of the data passed to the next, resulting in a streamlined and efficient process that optimally prepares input for the final classification using the Pizam-ANFIS model. The seamless interaction of these stages maximizes the accuracy and robustness of the overall system.

### 3.8. Action Unit Identification

Here, the CSE network determines the human AUs  $F_z^{sel}$  for quickly identifying emotions during training. Human action units encompass various expressions and movements

such as slit, eyes closed, squint, blink, wink, and others. They also include actions such as raising the inner and outer brows, lowering the brow, lifting the upper lid, wrinkling the nose, raising the cheeks, tightening the lids, and drooping the lids. A CNN is utilized for its capability to process high-dimensional data effectively without significant information loss. However, in existing CNNs, many neurons still need to be updated because the ReLU activation function does not preserve negative values due to its monotonic and linear nature. The proposed method utilizes Hard Swish (HS), which is nonmonotonic and smooth, to address this issue. The nonmonotonic property of HS stabilizes the network's gradient, allowing it to maintain small negative values. Additionally, the CNN's embedding and correlated interference modules are crucial for effective recognition. These enhancements ensure that the network can better capture and process the nuances of human action units, leading to more accurate and robust recognition. The correlated interference module received and processed the discriminative AU features' estimations from the embedded module. It calculates the correlations between the differentiating characteristics. As a result, the planned CNN is known as CSE.

#### 3.8.1. Input Layer

The input layer of a neural network consists of artificial neurons that introduce the initial data into the system, setting the foundation for processing by the subsequent layers of artificial neurons.

#### 3.8.2. Convolution Layer

In the convolution layer, an element-wise product is computed between each element of the kernel and the input array at every position within the tensor. The resulting products are then summed to produce the output value for the corresponding location in the output array. This process is repeated with multiple kernels to generate diverse feature maps. Then, convolution ( $L_{con}$ ) is expressed as,

$$L_{con} = \sum_u \sum_u (F_z^{sel})(g - u, h - u) * w(u, u) \quad (22)$$

Where  $g$  and  $h$  are the input matrix dimension size,  $w(u, u)$  represents the kernel having  $u \times u$  dimension size. The convolution parameters can reduce the model's complexity.

#### 3.8.3. Nonlinear Activation Function

The HS activation function is used for this purpose. The main task of using nonlinearity is to adjust or cut off the generated convolution output. The activation function is expressed in the mathematical representation as,

$$A = L_{con} \frac{R6(L_{con}+3)}{6} \quad (23)$$

Where A denotes the output of the HS activation function and R denotes the ReLU activation function.



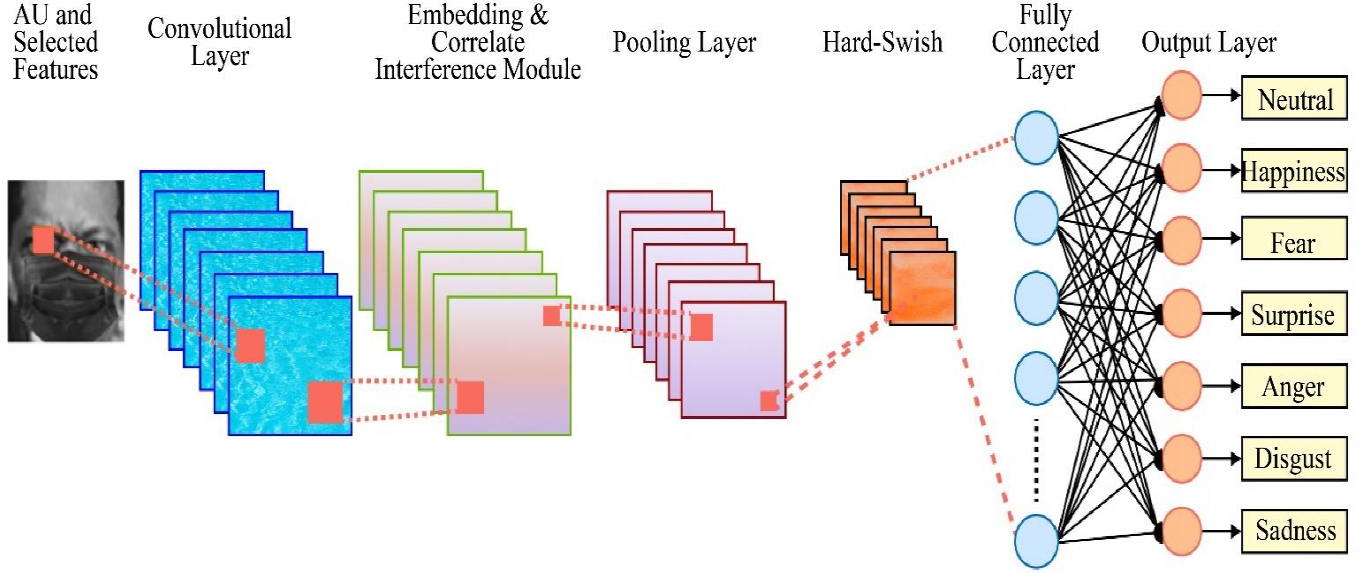


Fig. 2 Architecture of the proposed CSE

3.8.4. Embedding Module and Correlated Interference Module

In this step, the features derived from the nonlinear function are fed into the embedding module. This module utilizes a deeper convolutional network as a feature extraction mechanism, enhancing the capacity for feature representation by extracting discriminative AU features. Then, the output of the embedding module is calculated as,

$$A^{emb} = Embed\{(A) * w(u, u)\} \quad (24)$$

Here, Embed{ } signifies the embedding function. The discriminative AU features are given into the correlated interference module, which efficiently calculates the correlation between the features, and it is represented as,

$$A^{corr} = \frac{\sum x^{A^{corr}} y^{A^{corr}} - \sum x^{A^{corr}} \sum y^{A^{corr}}}{\sqrt{(\sum (x^{A^{corr}})^2 - (\sum x^{A^{corr}})^2)(\sum (y^{A^{corr}})^2 - (\sum y^{A^{corr}})^2)}} \quad (25)$$

Here,  $A^{corr}$  specifies the output of the correlated interference module.

3.8.5. Pooling Layer

The pooling layer executes a standard down-sampling action that lowers the in-plane dimensionality of the feature maps. It produces the highest value detected within the pooling filter, using this value as the result. The pooling ( $L_{pool}$ ) operation can be expressed as:

$$L_{pool} = \frac{A^{corr} - w}{s} + 1 \quad (26)$$

Where S represents the kernel strides, the process continues until it reaches the last layer.

3.8.6. Fully Connected Layer

The output feature maps from the final convolution or pooling layer are flattened into a one-dimensional array of numbers. The last completely linked layer has an equal number of output nodes corresponding to the number of classes. Calculating the flattened output as,

$$L_{fully} = L_{pool} - (w(u \times u) - 1) \quad (27)$$

Where,  $L_{fully}$  is the output of the fully connected layer.

3.8.7. Softmax Layer

The activation function, primarily used in the output layer, normalizes the real values in the range (0, 1) from the last fully connected layer into target class probabilities. This is achieved using the softmax function, which is defined by the following equation,

$$L_{soft} = \frac{e^{L_{fully}}}{\sum L_{fully}} \quad (28)$$

Where  $L_{soft}$  is the output of the softmax activation function. Later, the loss function is evaluated using the below equation,

$$loss = \|O^{target} - L_{soft}\| \quad (29)$$

Here,  $O^{target}$  specifies the target output. Finally, the identified AU is denoted as ( $L_{soft}$ ). The pseudocode of the proposed CSE is,

Input : Dimension-reduced features ( $F_z^{sel}$ )  
 Output : Action units ( $L_{soft}$ )



Begin

```

Initialize parameters  $L_{con}, w(u, u) L_{pool}$ 
Compute weight value
While  $j = 1toZ$ 
Forj = 1
    Compute convolution operation  $\eta$ 
        Evaluate activation function
         $A = L_{con} \frac{R6(L_{con}+3)}{6}$ 
        Compute Embedding Module
    Perform Correlated Interference Module
End for
While  $j = 2toZ$ 
Forj = 2
    Compute convolution operation  $\eta$ 
        Evaluate activation function
         $A = L_{con} \frac{R6(L_{con}+3)}{6}$ 
    Compute pooling operation  $L_{pool}$ 
End for
End while
Flattening all the layers
Evaluate the softmax activation function  $L_{soft}$ 
If ( $O^{t\ arg\ et} \neq O^{L_{soft}}$ )
    Stop criteria
}
else
{
    Set  $iter = iter + 1$ 
}
End if
Return  $L_{soft}$ 

```

End

### 3.9. Emotion Classification

The Pizam-ANFIS is used to categorize the types of emotions by taking the input as selected features and action units from the occluded and mask-covered input images once the action units have been identified. The Adaptive Neuro-Fuzzy Inference System (ANFIS) is a computational and predictive model that integrates the fuzzy Sugeno method with an adaptive neural network system. However, the adapted Sugeno fuzzy interference system introduces computational complexity while designing the higher-order fuzzy models.

The Mamdani fuzzy interference system in the defuzzification process is induced with modification in the existing ANFIS to avoid this issue. It uses the center of gravity technique for the defuzzification process, and the bell membership is replaced with the Piz membership function, which reduces the computational complexity and produces effective outcomes. Here, the second layer performs the fuzzification process, with the nodes in this layer being adaptive. The fuzzified output for the  $t^{th}$  layer  $\Phi_t$  is,

$$\Phi_t = \mu_1(\eta_{W_h}) \quad (30)$$

$$\Phi_t = \mu_2(\eta_{W_v}) \quad (31)$$

Where,  $\mu_1$  and  $\mu_2$  represent input node,  $W_h$  and  $W_v$  denotes the value of weights,  $\eta$  denotes the Piz membership function (layer1), and it is calculated as,

$$\eta = \frac{1}{1 + \left( \frac{AF^{points} - p1}{p2} \right)^2} \quad (32)$$

Here,  $AF^{points}$  denote the feature and AU points. In the third layer, the output signals from the previous layers are multiplied. This layer processes the outputs from the second layer  $\epsilon_t$ , resulting in:

$$\epsilon_t = \mu_1(\eta_{W_h}) * \mu_2(\eta_{W_v}) \quad (33)$$

The output of each node represents the firing strength of the rules. In the fourth layer, the output, described as the normalized firing strength  $(\epsilon_t)^*$ , is mathematically represented using the Radial Basis Function (RBF) as follows,

$$(\epsilon_t)^* = \sum_i \eta_{W_h} \zeta(\mu_2(\eta_{W_v}), \epsilon_t) + b \quad (34)$$

Here,  $\zeta$  and  $b$  denote kernel and bias. The consequent part of the fuzzy rules is executed in the fourth layer. The nodes in this layer are adaptive, and the node function is formulated as follows,

$$(\bar{\epsilon}_t)^* = \frac{(\epsilon_t)^*}{(\varphi_i \mu_1 + a_i \mu_2 + L_i)} \quad (35)$$

Where,  $\varphi_i$ ,  $a_i$  and  $L_i$  denote linear adaptive parameters,  $(\bar{\epsilon}_t)^*$  represent defuzzification using the Mamdani interference system's defuzzification process. Finally, the last layer predicts the emotions of the human ( $\Gamma$ ), and it is represented as,

$$\Gamma = \sum(\epsilon_t)^* (\varphi_i \mu_1 + a_i \mu_2 + L_i) \quad (36)$$

After training the proposed network, the image, which has to be tested, is given to the system for testing. By testing the data, the output layer classifies the emotions as Neutral, Happiness, Fear, Surprise, Anger, Disgust, and Sadness.

## 4. Results and Discussion

This section details the experiments performed on the PYTHON platform to validate the proposed scheme's performance. The experiments utilized a synthetic dataset created from publicly available sources. The dataset was divided into two parts: 80% of the images were used for training, while the remaining 20% were set aside for testing. Figure 3 depicts how sample images from the dataset were preprocessed and incorporated into the operation.

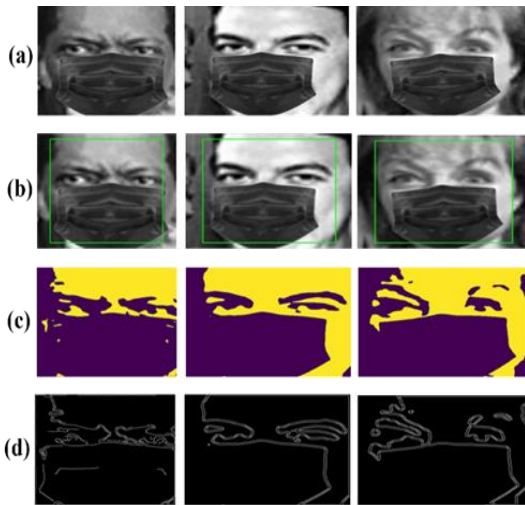


Fig. 3 Sample images of a human face with an emotion (a) Input images, (b) Face detected image, (c) Patch generated image, and (d) Edge detected image

4.1. Performance Analysis of Proposed CSE-Pizam-ANFIS

To thoroughly assess channel estimation performance, the anticipated CSE-Pizam-ANFIS algorithm was benchmarked against several well-established methods. These included the ANFIS, CNN, LSTM network, and ANN.

The efficacy and advantages of the CSE-Pizam-ANFIS approach in channel estimation were effectively validated by conducting a comprehensive comparison with these existing algorithms. Figure 4 presents a detailed assessment of the proposed CSE-Pizam-ANFIS model’s performance compared to existing models, focusing on key metrics like accuracy, precision, recall, sensitivity, specificity, and f-measure. Higher values in these metrics indicate more efficient model performance.

The proposed CSE-Pizam-ANFIS model achieves an impressive accuracy of 99.28%, which is notably higher than the accuracy of the existing models: ANFIS at 97.22%, CNN at 95.24%, LSTM at 93.34%, and ANN at 90.97%. In addition to accuracy, the proposed model excels in other metrics. It records a precision of 99.67%, a recall of 99.35%, a sensitivity of 99.35%, a specificity of 99.09%, and an f-measure of 99.51%.

These values surpass those of the existing models, showcasing the superior performance of the proposed model across all evaluated aspects. This comprehensive analysis underscores the effectiveness of the proposed model in AU classification and emotion classification tasks, significantly outperforming current alternatives.

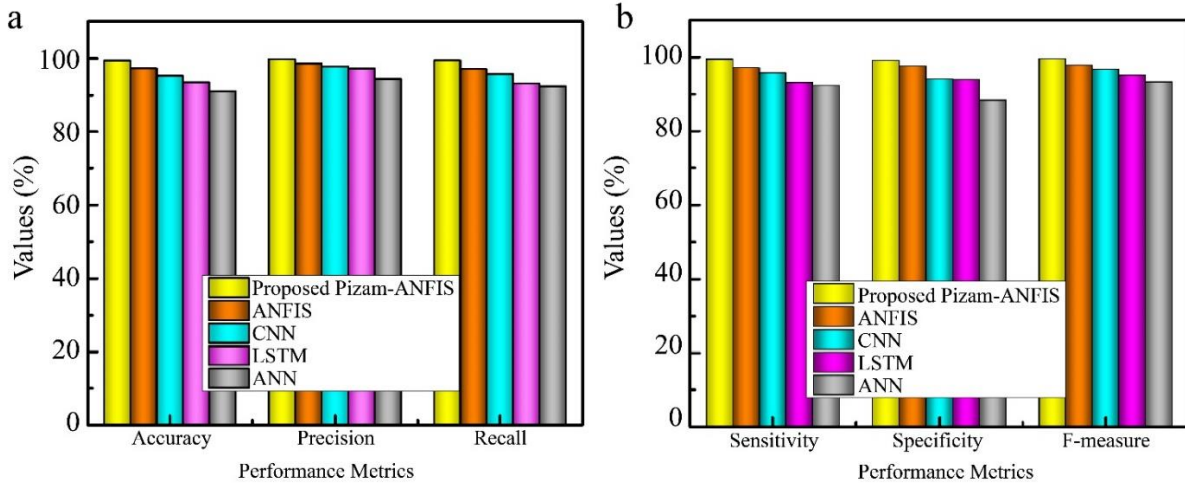


Fig. 4 Illustrative comparison of the proposed and existing models (right-hand side): (a) Accuracy, Precision, Recall metrics, and (b) Sensitivity, Specificity, and F-measure parameters

Table 1 presents a detailed performance evaluation of both the proposed and existing models using various metrics, including False Positive Rate (FPR), False Rejection Rate (FRR), False Negative Rate (FNR), Positive Predictive Value (PPV), Negative Predictive Value (NPV), and Matthews Correlation Coefficient (MCC). Higher values of FPR, FRR, and FNR indicate improved performance of the proposed model, while lower values of PPV, NPV, and MCC demonstrate its higher efficiency. For example, the proposed model shows a 63.78% improvement in FPR compared to ANFIS, 84.81% compared to CNN, and 92% compared to

ANN. Additionally, the FRR of the proposed model is 94.51% better than that of LSTM and other existing models. Similarly, the FNR, PPV, NPV, and MCC metrics for both the proposed and existing models have been analyzed and compared.

This detailed analysis reveals the superior efficiency and performance of the developed AU and emotion recognition system. Figure 5 illustrates the computational time analysis, comparing the proposed and existing models. Attaining a lower time for the proposed model indicates the efficient time of the proposed model.

Table 1. Performance evaluation of proposed and existing models

Techniques	FPR	FRR	FNR	PPV	NPV	MCC
Proposed CSE-Pizam-ANFIS	0.00900901	0.006451613	0.006452	0.996764	0.9821429	0.9871895
ANFIS	0.02484472	0.11764706	0.029412	0.985075	0.9515152	0.971719
CNN	0.05932203	0.04290429	0.042904	0.976431	0.8951613	0.8845861
LSTM	0.0610687	0.068965517	0.068966	0.971223	0.8601399	0.8544533
ANN	0.11764706	0.077192982	0.077193	0.942652	0.8450704	0.7963936

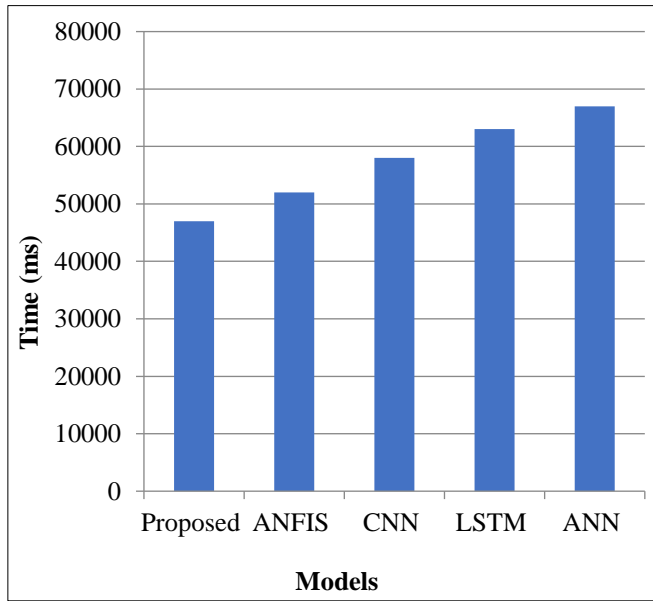


Fig. 5 Computational time analysis

The training time of the proposed model is given as 47015ms, whereas the existing ANFIS (52009ms), CNN (58006ms), LSTM (63006ms), and ANN (67010ms) take more time to train the proposed model. This can be achieved by inducing the HS and embedding a correlated interference module to stabilize the network’s gradient and efficiently recognize action units. Additionally, the Piz membership function and the Mamdani defuzzification method was introduced, which aids in the classification of emotions for computational complexity.

The confusion matrix in Figure 6 illustrates the high accuracy of the Pizam-ANFIS-based FER model, with the majority of predictions aligning correctly with actual labels. The strong diagonal values indicate precise classification across seven emotional categories-Angry, Disgust, Fear, Happy, Neutral, Sad, and Surprise-while minimal off-diagonal values reflect rare misclassifications.

The model effectively differentiates emotions even under partial occlusion and mask-wearing, proving its robustness. Compared to traditional approaches like CNN, LSTM, and ANN, Pizam-ANFIS demonstrates superior performance with 99.28% accuracy, making it a highly efficient FER solution for real-world applications.

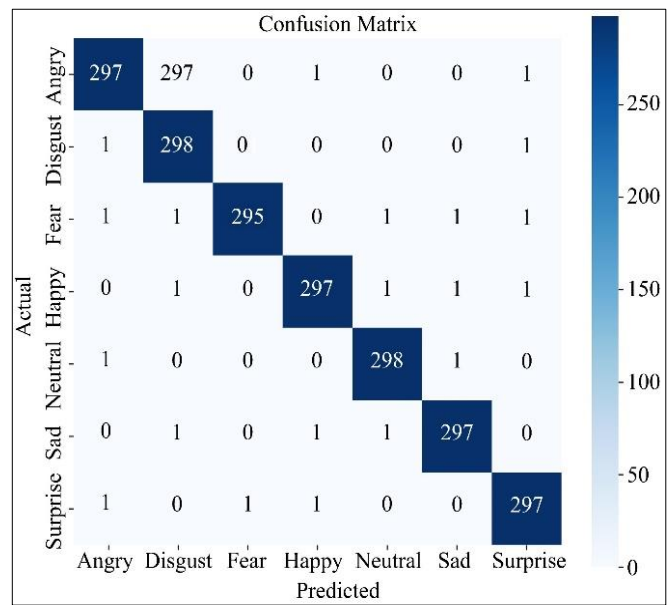


Fig. 6 Confusion matrix of the proposed Pizam-ANFIS model

4.2. Performance Analysis of Patch Generation

To highlight the advantages of the proposed model, the performance of the LSW-KMC was evaluated. Comparison was made with the existing models using metrics such as the Jaccard Index, Dice score, and Clustering Time. Table 2 presents a detailed analysis of the Jaccard Index for both the proposed and existing models.

The Jaccard Index measures the similarity between pixel groupings across different clusters, with one indicating that two clusters have perfectly extracted the same pixels and 0 indicating no overlap. The proposed model achieved a Jaccard Index of 0.03263298, demonstrating superior performance compared to the existing models, which showed lower coefficients. This result underscores the enhanced effectiveness of the new patch generation technique in accurately identifying similar clusters.

Table 2. Jaccard index

Method	Jaccard Index
Proposed LSW-KMC	0.03263298
K Means	-0.0690296
FCM	-0.068997
K Medoid	-0.0690228
CLARA	-0.069018

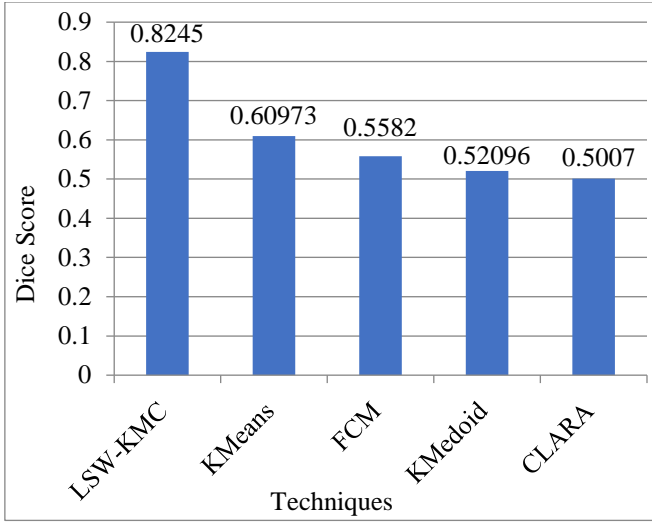


Fig. 7 Dice score analysis

Table 3. Clustering time analysis

Method	Clustering Time (ms)
Proposed LSW-KMC	38010
K Means	43005
FCM	47011
K Medoid	53010
CLARA	56012

Figure 7 comprehensively analyzes the Dice scores, comparing the proposed model with existing models. The Dice coefficient measures the pixel-wise agreement between predicted segmentation and the corresponding ground truth, showing that the proposed model achieved a Dice score of 0.8245. This performance significantly surpasses that of the existing models: K-Means with a score of 0.60973, Fuzzy c-Means (FCM) with 0.5582, K-Medoid with 0.52096, and Clustering Large Applications (CLARA) with 0.5007. This analysis highlights the superior performance of the proposed method. Furthermore, Table 3 presents the performance results, underscoring the proposed model’s efficiency in terms of clustering time. The proposed model’s clustering time is 38010ms, showing an improvement of 4995ms over K Means, 9001ms over FCM, and 18002ms over CLARA. This indicates that the LSW-KMC technique generates facial parts with greater accuracy and in a shorter time frame. The overall success of the proposed model is attributed to the careful selection and modification of existing patch-generation techniques, as established in previous studies. By refining these existing models, the proposed approach effectively generates accurate facial parts more efficiently.

**4.3. Comparative Evaluation of the Suggested and Earlier Approaches**

In this section, the performance of the proposed methodology with existing hybrid approaches developed by various researchers based on their classification accuracy is compared.

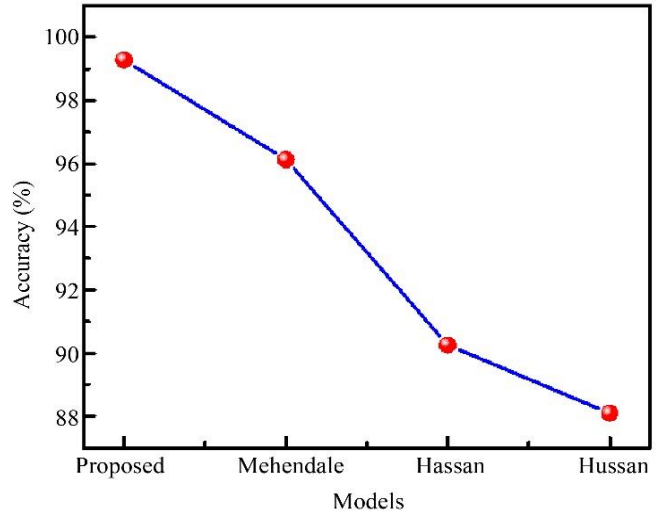


Fig. 8 Comparison of accuracy between proposed and existing models

Figure 8 illustrates the accuracy performance of the proposed framework under various conditions. The model consistently demonstrated superior performance across all tested scenarios.

The existing models utilized different techniques: Hussain & Salim Abdallah Al Balushi (2020) [15] employed a graph mining scheme, Hassan & Mohammed (2020) [14] used a CNN model, and Mehendale (2020) [11] applied FERC. In contrast, the proposed PAFEROFA model achieved higher accuracy in emotion classification, which is attributed to using CSE and PIGA techniques for recognizing AUs and selecting optimal features for training, respectively. Therefore, it is evident that the overall performance of the proposed methodology surpasses that of the existing techniques.

The proposed Pizam-ANFIS framework demonstrates superior performance compared to existing state-of-the-art techniques due to its innovative and robust methodology. Unlike traditional methods, which often struggle with occlusion and mask-related challenges, our framework integrates advanced techniques at every stage of the FER process. The LSW-KMC algorithm ensures precise segmentation by effectively isolating key facial regions, while the CGED method enhances edge detection, making it resilient to noisy and partially obscured inputs.

The PIGA-based dimensionality reduction technique also optimizes feature selection, retaining the most critical features while reducing computational overhead. When compared to models such as CNN, LSTM, ANN, and ANFIS, the Pizam-ANFIS framework consistently achieves higher performance metrics. For example, it achieves an impressive accuracy of 99.28%, surpassing CNN’s 95.24% and ANFIS’s 97.22%. Furthermore, it records lower false positive and false rejection rates, reflecting its enhanced reliability across diverse scenarios. Including the Pizam-ANFIS classifier, featuring

optimized defuzzification processes and the Piz membership function, is crucial in improving classification accuracy while maintaining computational efficiency. These findings underscore the effectiveness of the proposed approach in overcoming the limitations of existing models, demonstrating its potential for robust and accurate emotion recognition in real-world applications.

#### 4.4. Ethical Implications of FER Technology

While the technical advancements in occlusion handling significantly contribute to the field of FER, it is equally important to consider the ethical implications of emotion recognition technology. Deploying FER systems raises concerns about privacy, consent, and potential misuse. For instance, the indiscriminate use of emotion recognition in surveillance without clear consent could infringe on individuals' privacy rights. Additionally, misinterpreting emotions due to algorithmic biases could lead to unjust outcomes in sensitive contexts such as law enforcement or recruitment. Ethical questions also arise around the storage and handling of facial data, especially when it pertains to minors or vulnerable populations. To address these concerns, FER implementations must prioritize transparency, seek informed consent, and comply with data protection regulations such as GDPR. Moreover, ensuring that these systems are designed with fairness and inclusivity in mind is crucial, reducing biases linked to demographic factors like age, gender, or ethnicity. Balancing innovation with ethical considerations will foster trust and pave the way for responsible deployment of FER technologies.

### 5. Conclusion

FER is a crucial method for assessing emotional states. However, traditional recognition models often struggle with accuracy due to challenges like partial occlusion and wearing

face masks. To address these issues, a novel FER method is developed. The process begins with preprocessing the input image and detecting the face. Differential parts of the face are then extracted using the LSW-KMC method, which identifies a similarity score of 0.0326 within a time frame of 38010ms. Following this, feature extraction and selection are performed using the PIGA technique, which is known for its high efficiency.

These selected features classify Action Units (AUs) with a trained neural network model. Later, the features and AUs are fed into the Pizam-ANFIS classifier to determine the emotions. The proposed CSE-Pizam-ANFIS model achieved an impressive accuracy of 99.28% and a computation time of 47015ms. The proposed FER system demonstrated high efficacy, even under the challenging conditions of partial occlusion and face masks. Therefore, the proposed model outperforms existing methods. Currently, the model is designed to recognize emotions in individual subjects. Future research will focus on advancing emotion recognition from video inputs involving multiple people.

#### Availability of Data and Material

The data utilized in this study can be accessed upon request from the corresponding author.

#### Author's Contribution

JSB collected publicly available datasets for emotion recognition (e.g., happy, sad, angry) and developed Python methods to evaluate the model's performance using metrics such as accuracy, precision, and recall for each emotion category. Additionally, JSB authored the complete manuscript and addressed the reviewers' comments. PLP supervised the experiments, reviewed and provided feedback on manuscript drafts, and guided the submission process.

### Reference

- [1] Wafa Mellouk, and Wahida Handouzi, "Facial Emotion Recognition Using Deep Learning: Review and Insights," *Procedia Computer Science*, vol. 175, pp. 689-694, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [2] Rohan Appasaheb Borgalli, and Sunil Surve, "Deep Learning for Facial Emotion Recognition Using Custom CNN Architecture," *Journal of Physics: Conference Series, 2<sup>nd</sup> International Conference on Computational Intelligence & IoT*, vol. 2236, pp. 1-12, 2022. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [3] Shrey Modi, and Mohammed Husain Bohara, "Facial Emotion Recognition Using Convolution Neural Network," *Proceedings - 5<sup>th</sup> International Conference on Intelligent Computing and Control Systems*, Madurai, India, pp. 1339-1344, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [4] Erick C. Valverde, Miracle Udurume, and Wansu Lim, "Performance Analysis of a Deep Learning-Based Facial Emotion Recognition System on Edge Device," *Neural Computing and Applications*, pp. 695-696, 2022. [[Google Scholar](#)] [[Publisher Link](#)]
- [5] Shan Li, and Weihong Deng, "Deep Facial Expression Recognition: A Survey," *IEEE Transactions on Affective Computing*, vol. 13, no. 3, pp. 1195-1215, 2022. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [6] Anjad Rehman Khan, "Facial Emotion Recognition Using Conventional Machine Learning and Deep Learning Methods: Current Achievements, Analysis and Remaining Challenges," *Information*, vol. 13, no. 6, pp. 1-17, 2022. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [7] Mahmut Dirik, "Optimized Anfis Model with Hybrid Metaheuristic Algorithms for Facial Emotion Recognition," *International Journal of Fuzzy Systems*, vol. 25, pp. 485-496, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]



- [8] Shervin Minaee, Mehdi Minaei, and Amirali Abdolrashidi, "Deep-Emotion: Facial Expression Recognition Using the Attentional Convolutional Network," *Sensors*, vol. 21, no. 9, pp. 1-16, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [9] Elham S. Salama et al., "A 3D-Convolutional Neural Network Framework with Ensemble Learning Techniques for Multi-Modal Emotion Recognition," *Egyptian Informatics Journal*, vol. 22, pp. 167-176, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [10] Fakir Mashuque Alamgir, and Md. Shafiu Alam, "An Artificial Intelligence-Driven Facial Emotion Recognition System Using Hybrid Deep Belief Rain Optimization," *Multimedia Tools and Applications*, vol. 82, pp. 2437-2464, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [11] Ninad Mehendale, "Facial Emotion Recognition Using Convolutional Neural Networks (FERC)," *SN Applied Sciences*, vol. 2, no. 3, pp. 1-8, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [12] Xiao Liu, Xiangyi Cheng and Kiju Lee, "GA-SVM-Based Facial Emotion Recognition Using Facial Geometric Features," *IEEE Sensors Journal*, vol. 21, no. 10, pp. 11532-11542, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [13] Abdulrahman Alreshidi, and Mohib Ullah, "Facial Emotion Recognition Using Hybrid Features," *Informatics*, vol. 7, no. 1, pp. 1-13, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [14] Alia K. Hassan, and Suhaila N. Mohammed, "A Novel Facial Emotion Recognition Scheme Based on Graph Mining," *Defence Technology*, vol. 16, no. 5, pp. 1062-1072, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [15] Shaik Asif Hussain, and Ahlam Salim Abdallah Al Balushi, "A Real-Time Face Emotion Classification and Recognition Using Deep Learning Model," *Journal of Physics: Conference Series, First International Conference on Emerging Electrical Energy, Electronics and Computing Technologies*, Melaka, Malaysia, vol. 1432, no. 1, pp. 1-14, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [16] Bitu Houshmand, and Naimul Mefraz Khan, "Facial Expression Recognition Under Partial Occlusion from Virtual Reality Headsets Based on Transfer Learning," *Proceedings-IEEE 6<sup>th</sup> International Conference on Multimedia Big Data*, New Delhi, India, pp. 70-75, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [17] G.S. Monisha et al., "Enhanced Automatic Recognition of Human Emotions Using Machine Learning Techniques," *Procedia Computer Science*, vol. 218, pp. 375-382, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [18] M.A.H. Akhand et al., "Facial Emotion Recognition Using Transfer Learning in The Deep CNN," *Electronics*, vol. 10, no. 9, pp. 1-19, 2021. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [19] Soumyajit Saha et al., "Feature Selection for Facial Emotion Recognition Using Cosine Similarity-Based Harmony Search Algorithm," *Applied Sciences*, vol. 10, no. 8, pp. 1-22, 2020. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [20] Chahak Gautam, and K.R. Seeja, "Facial Emotion Recognition Using Handcrafted Features and CNN," *Procedia Computer Science*, vol. 218, pp. 1295-1303, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [21] Giovanna Castellano, Berardina De Carolis, and Nicola Macchiarulo, "Automatic Facial Emotion Recognition at the COVID-19 Pandemic Time," *Multimedia Tools and Applications*, vol. 82, pp. 12751-12769, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [22] Shrouk Wally et al., "Occlusion Aware Student Emotion Recognition Based on Facial Action Unit Detection," *arXiv*, pp. 1-14, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [23] Yasmeen ELSayed, Ashraf ELSayed, and Mohamed A. Abdou, "An Automatic Improved Facial Expression Recognition for Masked Faces," *Neural Computing and Applications*, vol. 35, no. 20, pp. 14963-14972, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [24] Mukhriddin Mukhiddinov et al., "Masked Face Emotion Recognition Based on Facial Landmarks," *Sensors*, vol. 23, no. 3, pp. 1-23, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]
- [25] Xiaoliang Zhu et al., "Hybrid Domain Consistency Constraints-Based Deep Neural Network for Facial Expression Recognition," *Sensors*, vol. 23, no. 11, pp. 1-16, 2023. [[CrossRef](#)] [[Google Scholar](#)] [[Publisher Link](#)]