# Multi-Featured Extraction and Convolution Neural Network Based SVM for Automatic Facial Emotion Recognition

## M.Regina[1*], M.S. Josephine[2], V. Jeyabalraja[3]

[1,2]Dept. of computer Application, Dr.M.G.R Education and research Institute University, Chennai, India
[3]Dept. of computer science and Engineering, Velammal Engineering College, Chennai, India

*Corresponding Author: reginamathew@loyolacollege.edu*

*Abstract*—Human emotions are mental states of feelings that are exposed unconsciously and followed by physical changes in their facial muscles which entail the expressions on the face. Certain emotions commonly expressed by human are happiness, sadness, anger, fear, disgust, surprise, and neutral. For a non-verbal communication, facial expression plays a vital role since it appears because of inner core feelings of a person that reflects on the faces. For the automatic recognition of facial emotions, many methods are used such as Artificial Neural networks, Neuro-fuzzy, Wavelet transformation, etc. However, the existing methods take more time for data classification, low accuracy in the optimization process and high level of error rate. To overcome these concerns, this paper depicts an amphibious operation of Multi Support Vector Machine (SVM) with the Convolutional Neural Networks (CNN). Initially, the characteristics of the pre-processed face image are efficiently extracted by using Local Binary Pattern (LBP), Principal Component Analysis (PCA) and Gray Level Occurrence Matrix (GLCM). In this model, CNN works as a trainable feature extractor, and Multi-SVM performs as a recognizer. The proposed system's performance is analyzed with various human faces using the MATLAB tool. The results prove that the proposed method surpasses the earlier methods regarding high accuracy with low computation time and low error rate.

*Keywords*—*Convolutional Neural Networks, Face recognition, Local Binary Pattern, Principal Component Analysis.*

## I. INTRODUCTION

The investigation of facial expressions had its first noteworthy outcome with Darwin's exploration [1] on its effect in the advancement of the species as a kind of nonverbal communication. This communicating through recognizing a feeling in the facial expression, significantly quicker than verbal communication, would convey an extraordinary preferred standpoint to the human species, guessed Darwin. Under this evaluative viewpoint, the facial expressions would be all-inclusive to all people. Late investigations [2] support this hypothesis, with researchers recognizing six all-inclusive facial expressions: Happiness, Sadness, Fear, Surprise, Anger, and Aversion. Enthusiastic discernment is a critical component of human communication, used to translate occasions, social connections and to human relations generally [3].

Emotional expression has a few activities related with it, for example, face actions and body motions, variations in voice tone, physiological changes in the skin obstruction and facial flushing, to refer to a few. Emotional perception is an unpredictable activity that may include a few components. Constraining the investigation of emotional expressions to the perception of the condition of facial muscles [4],

emotional perception understanding can be confined to the examination of consecutive images or even single images.
Other than the strategies utilizing deep architecture, there are numerous others in the literature. However, a few parts of the assessment of these techniques still merit consideration. [5] In this specific circumstance, systems for expression recognition are normally founded on appraisals of the developments of the facial muscles [6] and the eyes [7], or on measures to build up a connection between the state of parts of the face and the emotions. This data can be acquired through still images or through arrangements of images that demonstrate the feeling going from neural to its apex.

As of late, computational intelligence strategies, for example, deep neural networks (DNN) [8], specifically convolutional neural networks (CNNs), have been utilized to extricate includes more successfully than physically planned ad hoc extractors. The architecture of CNNs requires a broad arrangement of parameters that are found out from an expansive arrangement of already named data, for this situation, an arrangement of beforehand marked images. In a few circumstances, this can be an issue, since it might be hard to discover public datasets with a gigantic measure of images. A system to overcome this constraint is the artificial enlargement of the database via label-preserving

modifications in the data. This procedure is regularly known as data augmentation. [9]

In this paper, an idea is provided for the recognition of human emotions through facial expressions (Frames by frames) using a Convolution Neural Network (CNN) and the process of human thought variation. This work frames ensemble of classifiers is utilized for classifying the various emotions from the video sequence. The Gabor wavelet is used for the temporal feature (Gabor features) extraction of eye and mouth, and z-score normalization is applied to generate the feature vector. Finally, from the obtained Gabor features the universal emotions are classified accordingly.

The organization of this paper is as follows, Section I illustrates the introduction, Section II describesabout the literature review, Section III explains the methodology, Section IV gives the results and discussion, SectionV describes the conclusion and future work.

## II. LITERATURE REVIEW

Reference [10] exhibited a framework for perceiving feelings through facial expressions showed in live video streams and video sequences. The framework depended on the Piecewise B'ezier Volume Deformation tracker and was reached out with a Haar face detector to at first find the human face naturally. Their trials includes Naive Bayes and the Tree-Augmented-Naive Bayes (TAN) classifiers in person-dependent and person-independent tests on the Cohn-Kanade database demonstrated that great classification results could be acquired for facial emotion recognition.

Reference [11] concentrated on automatic facial expressions recognition from live video input utilizing temporal cues. Strategies for utilizing temporal information was widely investigated for discourse recognition applications. These techniques were template coordinating utilizing dynamic programming strategies and concealed Markov models (HMM). The work abused existing techniques and proposed another engineering of HMMs for automatically segmenting and perceiving human facial expressions from video successions. They investigated individual ward and individual free recognition of expressions.

Reference [12] displayed results on an independent completely automatic framework for continuous acknowledgment of facial activities from the Facial Action Coding System (FACS). The framework automatically identified frontal faces in the video stream and codes every frame regarding 20 Action units. Fundamental outcomes were introduced on an assignment of facial activity discovery in unconstrained articulations amid talk. Support vector machines and AdaBoost classifiers were looked at. For the two classifiers, the yield edge anticipated activity unit intensity.

Reference [13] showed that LBP highlights were compelling and proficient for facial expression acknowledgment. They additionally detailed Boosted-LBP to extricate the most discriminant LBP highlights and the best recognition execution was acquired by utilizing Support Vector Machine classifiers including Boosted-LBP highlights. Additionally, they explored LBP highlights for low-resolution facial expression recognition, or, in other words, issue yet only from time to time tended to in the current work. It was seen in the tests that LBP highlights perform steadily and powerfully over a helpful scope of low resolutions of face pictures, and yield promising execution in compacted low-resolution video sequences caught in genuine conditions.

Reference [14] presented Deep-learning-based FER approaches which exceptionally decreased the dependence on face-physics-based models and other pre-handling systems by empowering "end-to-end" figuring out how to happen in the pipeline straightforwardly from the input pictures.

Reference [15] focussed on accomplishing great accuracy while requiring just a little example data for training. Scale Invariant Feature Transform (SIFT) highlights were utilized to expand the execution on little data as SIFT does not require broad training data to create valuable highlights. In this paper, both Dense SIFT and regular SIFT were contemplated and contrasted when blended and CNN highlights.

Reference [16] proposed a two-part network comprising of a DNN-based design pursued by a Conditional Random Field (CRF) module for facial expression recognition in recordings. The trial results demonstrated that falling the profound network engineering with the CRF module extensively expanded the recognition of facial expressions in recordings and in particular it outflanks the best in class strategies in the cross-database tests and yielded equivalent outcomes in the subject-free analyses.

## III. METHODOLOGY

This paper presents a hybrid model of feature extraction using Local Binary Pattern (LBP),PCA and Gray Level Occurrence Matrix (GLCM). Firstly, the face image is pre-processed. In pre-processing, the image is resized, and gray conversion is done for further process. Every image has noise. It should be reduced for better prediction. Hence bilateral filter is used in the image to reduce noise. Then pre-processed image extracts features using Local Binary Pattern (LBP), PCA and Gray Level Occurrence Matrix (GLCM) for efficient feature extraction values. The extracted feature is used for the classification. Two superior classifiers: Convolutional Neural Network (CNN) and Multilevel Support Vector Machine (Multi-SVM), which have proven results in recognizing different types of patterns are used. In

this model, CNN works as a trainable feature extractor, and Multi-SVM performs as a recognizer. The new feature from the raw images are automatically extracted and predictions are generated by this hybrid model.
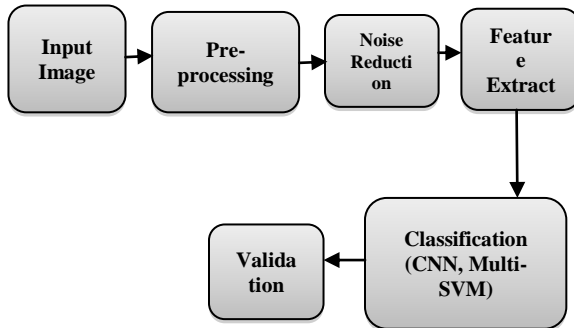


Fig. 1. Block Diagram of Proposed Methodology

The steps involved in the proposed methodology are as follows:
**Step1:** Collecting the images related to face emotions.
**Step2:** Pre-processing the input facial emotion images.
**Step3:** Reduction of noise using Bilateral Filter.
**Step 4:** Feature extraction using PCA, GLCM, and LBP techniques.
**Step5:** Classification using CNN, and Multi-SVM.
**Step 6:** Output Validation**.**

*A. Pre-processing:*
In pre-processing stage, input images are rehanged in their size and shape. The reshaped images are in the 256 x 256 range.In the pre-processing stage the image was resized and converts to gray if the image was in color. In this stage, image smoothing is done by noise removal, test removal and shaping of the image. The bilateral filter is used to remove noise in the image.

*1) Enhancement*
Enhancement is performed to make images more clearly and for excellent presentation of all digital processing in computers. In this stage, the region of the face with emotions are detected and cropped. Thereby, it becomes much useful and important for recognition of facial expressions. Hence, enhancement improves the quality of a given image. This process can be done by removing noise, emphasizing edges, modifying shapes, and enhancing contrast.

*2) Smoothing*
Smoothing is finished by noise removal, test removal and shaping of the image. The image is resized and changed over to gray if the image was in color. Bilateral filter is utilized to evacuate noise in the image. It is a non-clear, edge-saving, and noise decreasing smoothing filter for images. It substitutes the intensity of every pixel with a weighted average of intensity esteems from nearby pixels. This weight

relies on a Gaussian conveyance. The weights rely on the radiometric contrasts (e.g., range differences, similar to depth separate, color intensity, and so on.) and on Euclidean separation of pixels. In this manner, it holds sharp edges. It joins colors or gray levels relying upon both photometric comparability and geometric closeness and, and favors approximate values in both range and space.
The bilateral filter is defined as

$$I^{filtered}(x) = \frac{1}{W_p} \sum_{x_i \in \Omega} I(x_i) f_r(\| I(x_i) - I(x) \|) g_s(\| x_i - x \|)$$

(1)

Where the normalization term

$$W_p = \sum_{x_i \in \Omega} f_r(\| I(x_i) - I(x) \|) g_s(\| x_i - x \|) \ (2)$$

Equation (2) ensures that the filter preserves the image energy.
Where,
$I^{filtered}$ - image that is filtered;

$I$ - actual input image to be filtered;
$x$ - coordinates of the current pixel to be filtered;
$\Omega$ - centered window;

$f_r$ - range kernel for intensities' smoothing differences (usually Gaussian function);

$g_s$ - spatial kernel for coordinates' smoothing differences (usually Gaussian function).
The distortion-free clear image with the high contrast and quality is the outcome of this phase.

*B. Feature extraction*
Feature extraction means to decrease the original data by computing particular properties that separate one pattern from another pattern. In this stage, the feature is extracted utilizing the half and half of PCA, LBP, and GLCM.
PCA method is as often possibly utilized in signal processing to the data dimension minimization or data de-correlation.In signal processing, it can be defined as a transform of a given set of $n$ input vectors (variables) with the same length $K$ formed in the $n$-dimensional vector $x = [x_1, x_2, \ldots x_n]^T$ into a vector $y$ according to

$$y = A(x - m_x) \quad (3)$$

Every row of the vector $x$ consists of $K$ values corresponding to one input. The vector $m_x$ in Eq. (3) is the vector of average values of all input variables defined by the relation

$$m_x = E\{x\} = \frac{1}{k}\sum_{k=1}^{K} x_k$$
(4)

The GLCM system is the extraction of second-order statistical texture characteristics. GLCM displayed by Haralick incorporates data about the course of action of pixels having indistinguishable gray level qualities. In GLCM calculation, all the pixel pairs are counted, in that I esteem is relegated to the first pixel, and its matching pair has appointed by estimation of j which is replaced from the principal pixel by d. The included pair is set in the network Pd [i, j] of ith push and jth segment. As there is no important to contain the same number of pixels pairs in gray levels of [i, j] and [j, i], Pd [i, j] is unbalanced. The characteristics of Pd [i, j] can be institutionalized by isolating each passage by the whole number of pixel pairs present. Normalized GLCM N [i, j] is denoted by:

$$N[i, j] = \frac{P[i, j]}{\sum_{i}\sum_{j} P[i, j]}$$
(5)

LBP works with the 3×3 neighborhood. The pixel values of 8 neighbors are thresholded by the estimation of the middle pixel, at that point, the so-thresholded binary values are weighted by powers of two and consolidated to attain the LBP code of the middle pixel. Let $g_c$ and $g_{0,\ldots,}g_7$ denote individually the gray values of the middle and its eight neighbor pixels, at that point the LBP code for the middle pixel with coordinate $(x, y)$ is computed by

The pixel values of eight neighbors are thresholded by the estimation of the middle pixel, at that point, the so-thresholded binary values are weighted by forces of two and consolidated to attain the LBP code of the middle pixel. Let $g_c$ and $g_{0,\ldots,}g_7$ mean individually the gray values of the middle and its eight neighbor pixels, at that point the LBP code for the middle pixel with coordinate $(x, y)$ is computed by

$$LBP(x, y) = \sum_{p=0}^{7} s(g_c - g_p)2^p$$
(6)

where $s(z)$ - threshold function

$$s(z) = \begin{cases} 1, z \geq 0 \\ 0, z < 0 \end{cases}$$
(7)

In conventional real tasks, the statistic depiction of LBP codes, LBP histogram (LBPH), usually applied. The LBP codes of entire pixels to an input image are composed into a histogram as a texture descriptor, i.e.,

$$LBPH(i) = \sum_{x,y} \delta\{i, LBP(x, y)\}, i = 0,\ldots,2^7$$
(8)

Where $\delta(.)$ is the kroneck product function. [17]

*C. Classification*

Classification increases the accuracy, and better information from the individual class is achieved by using textures. Classification results rely upon several classes chosen by the user. Hence to have a good classification quality, it's important to choose an exact number of classes K. In this work, the deep conventional neural network is used for classification of the image. CNNs utilize a change of multilayer perceptrons designed for the need of minimal preprocessing. They are otherwise called space invariant artificial neural networks (SIANN) or shift invariant, depending on their architecture of shared-weights and translation invariance characteristics. Thereexists three types of CNNs: Pooling, Convolutional, and Fully-concentrated layers. In the convolution layer, the feature maps of previous layers' are convolved with learnable kernels and put via the activation function to produce output feature map.

Each result map might combine with the convolutions of many input maps. Generally, it is given by

$$x_j^l = f(\sum_{i \in M_j} x_i^{l-1} * k_{ij}^l + b_j^l$$
(9)

Where $M_j$ is the selection of input maps. MCSVM problem is to build a decision function given N samples especially with noise: $(x_1, y_1),\ldots,(x_N, y_n)$ where $x_i : i = 1,\ldots,N$ is a vector of length $n$ and $y_i \in \{1,\ldots,M\}$ denotes the class of the sample. The classical approach to examine MCSVM classification issues is to take into account the issue as a collection of binary classification issues. In the OAA method, one builds $M$ classifiers, one for each class. The $m^{th}$ classifier builds a hyperplane between class m and the $M-1$ remaining classes. [18]

## IV. RESULTS AND DISCUSSION

Input images to the proposed architecture comprise of grayscale cropped images with a width of 156 pixels and height of 176 pixels taken from the database, which contains faces communicating the feelings "Sad," and "Normal" Image Cropping is important to restrict all of them to the area around the eyes, mouth, and nose of the subject (Fig. 2). Automatic cropped algorithm dependent on Haar cascade face detection was produced with the end goal to encourage this assignment.

a                                b                                c
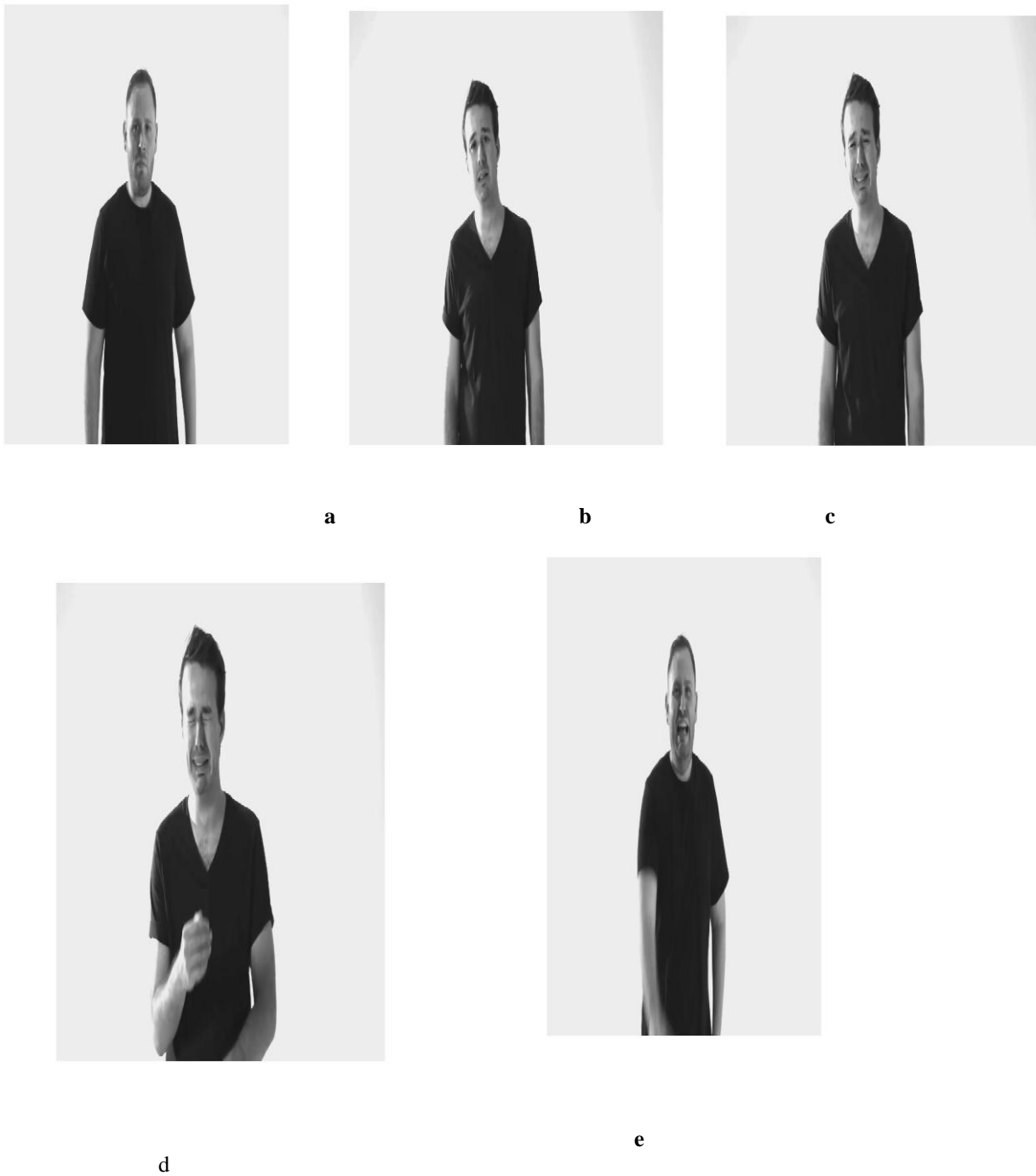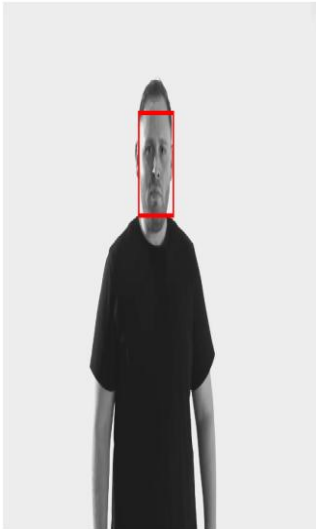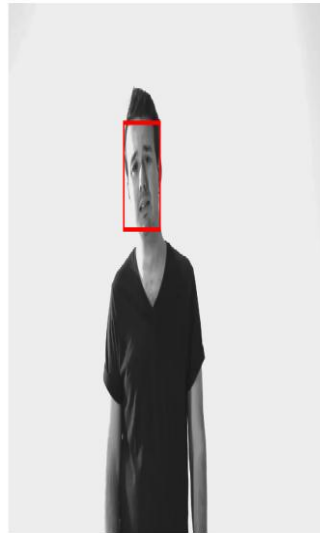


e

d

Fig 2. Input images

Emotions must be identified in face so that the region of the face showing emotions are cropped. Cropping of the face in images is done as shown in figure 3. Further the eyes, nose and mouth are converted into points which leads to the recognition of emotion.
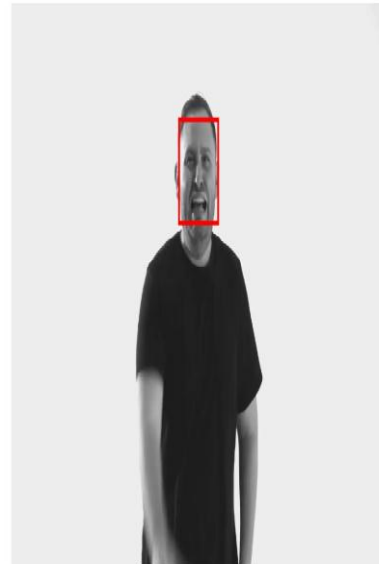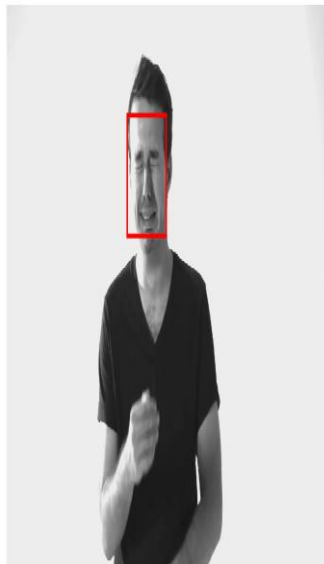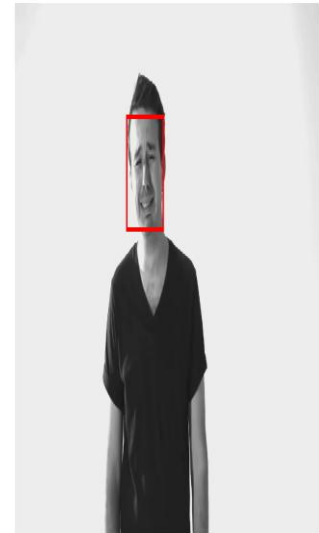
a           b           c



d           e

Fig 3. Cropping of the face in images

The cropped imagesshowing eyes, nose and mouth alone are given in figure 4. The cropped image are then applied for validation. The nose, eyes and mouth are recognized first and then the results are validated.
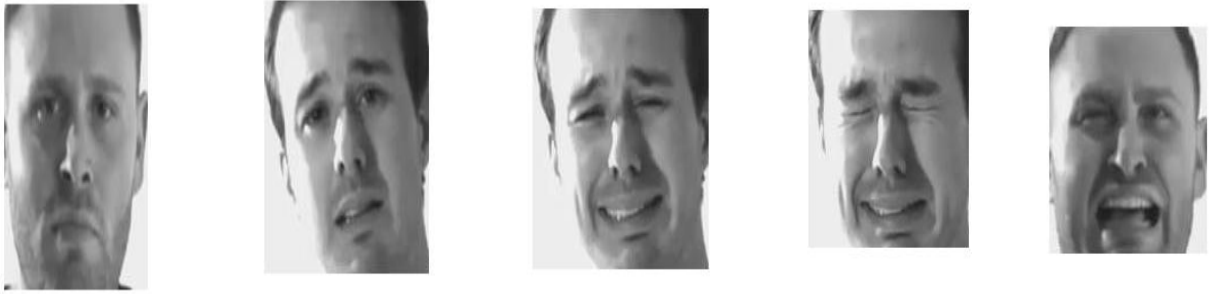
Fig 4. Cropped image of the faces

The results of all the cropped imagesare displayed as in figure 5 shown below indicating the corresponding face emotions.
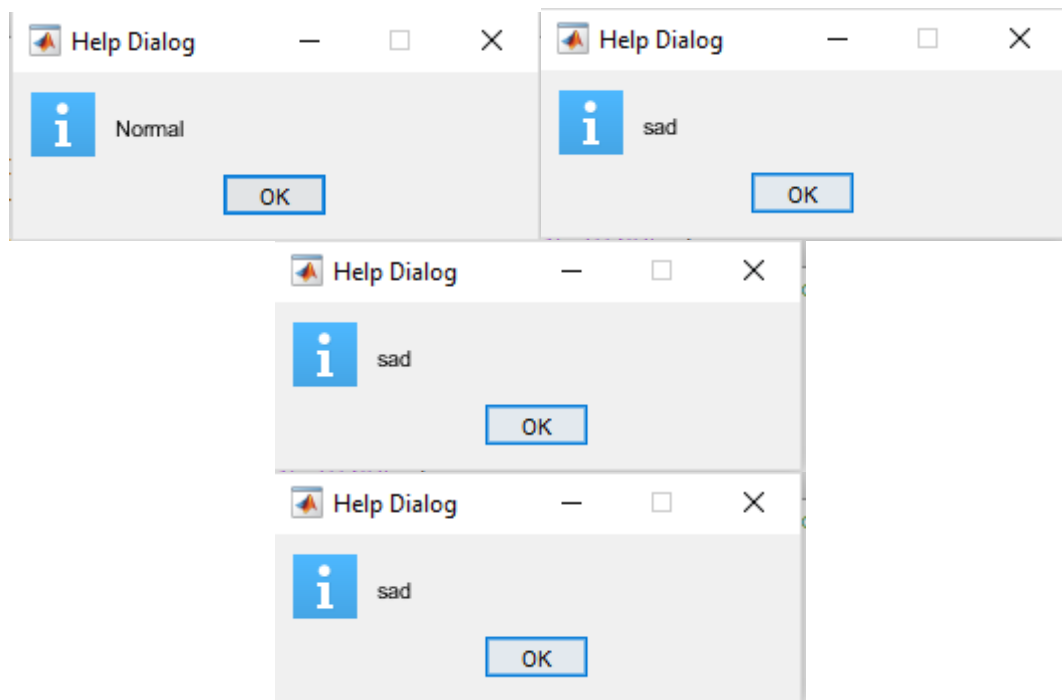


Fig 5. Results for images

**Table.1 Comparison of accuracy in existing and proposed method**

| Images | Existing method | Proposed Method |
|--------|-----------------|-----------------|
| a | 79% | 84% |
| b | 89% | 91% |
| c | 83% | 90% |
| d | 82% | 92% |
| e | 75% | 89% |

From the comparison table, it is understood that the accuracy rate is higher in all images of the proposed system than the existing method.
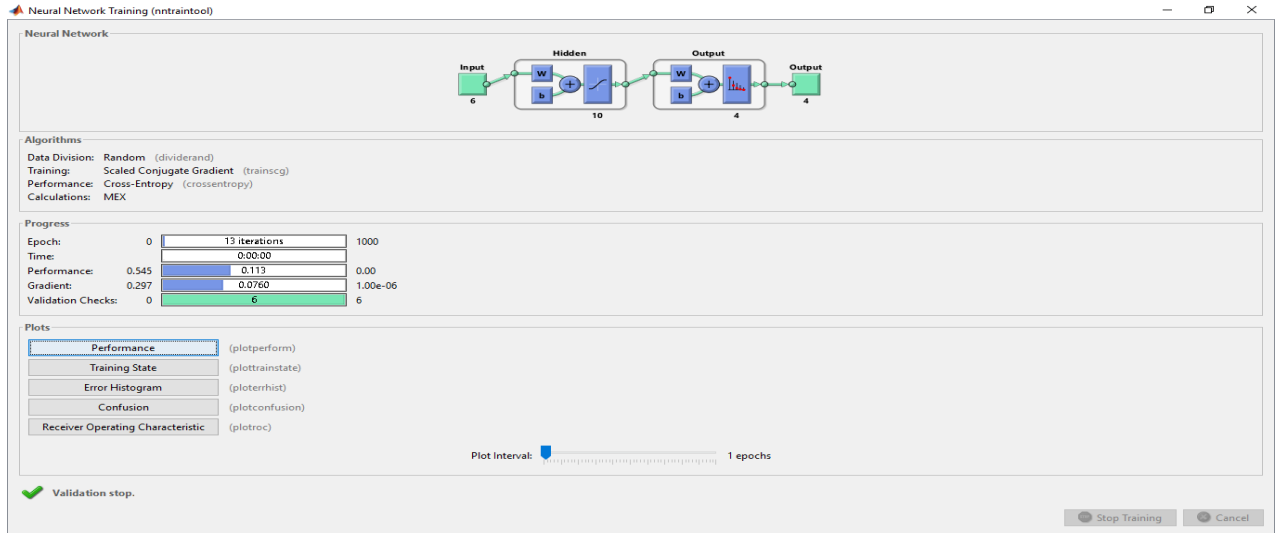
**Figure 7. Convolutional Neural Network**

When using all 16 attributes and all 4 classes, the confusion matrix obtained from the k-nearest neighbors clustering is shown in Figure 3, where along the x-axis are registered the true class labels and besides that, the y-axis are the k nearest .

neighbors class predictions. Along the first diagonal are the correct classifications, whereas all the other entries show misclassifications. The bottom right cell shows the overall accuracy.
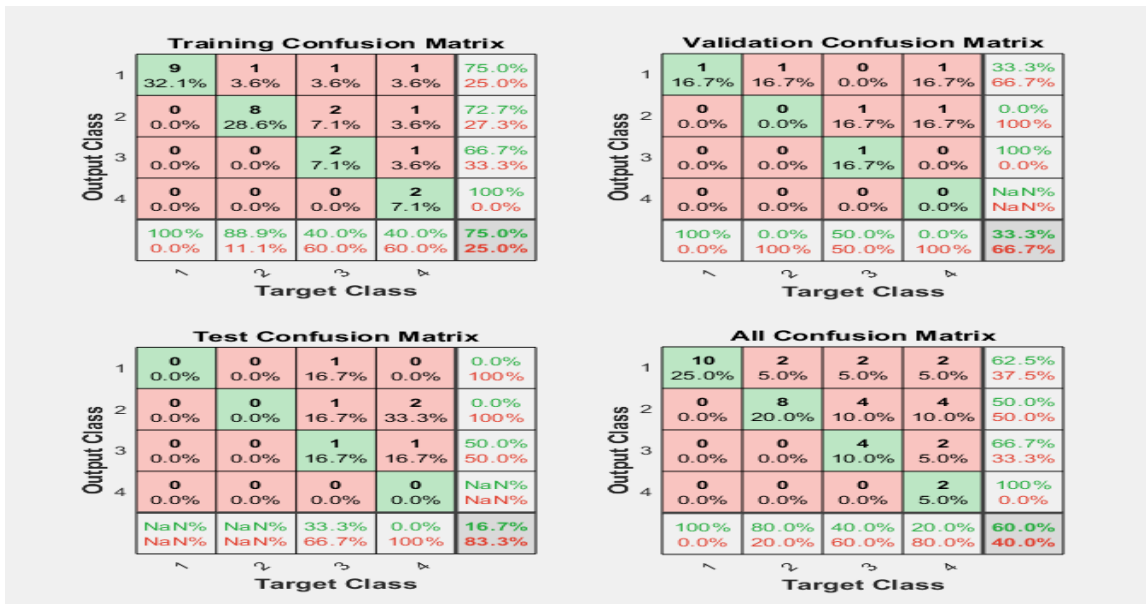


**Figure 8. Confusion Matrix**

**Table 2 Representation of Confusion Matrix**

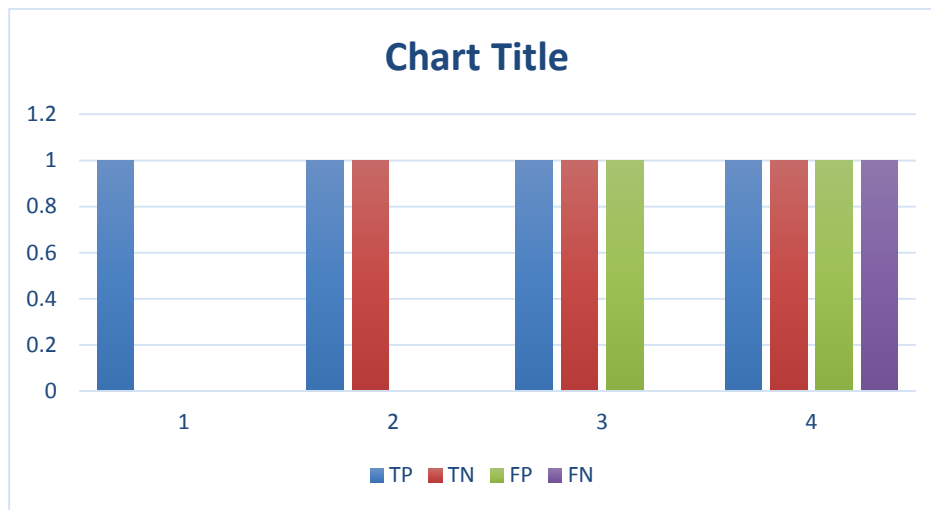| TP | TN | FP | FN | Condition |
|----|----|----|----|-----------|
| 1 | 1 | 1 | 1 | Normal |
| 0 | 1 | 1 | 1 | Casual |
|  |  |  |  |  |
| 0 | 0 | 1 | 1 | Sorrowful |
| 0 | 0 | 0 | 1 | Sad |

Fig 9. Graph for Confusion Matrix

From the confusion matrix, it is observed that the true positive, true negative, false positive, and false negative conditions are obtained in percentages. For True positive, the possibility is obtained as 37.5%, for True negative, it is obtained as 50%, for false positive, it is obtained as 33.3% and for false negative it is 0%. From these values, the facial emotions are identified.

## V. CONCLUSION

Human emotions like happiness, sadness, anger, fear, disgust, surprise, and neutral are being recognized. This is done by various methods such as Artificial Neural networks, Neuro-fuzzy, Wavelet transformation, etc. Thereby reducing the time required for data classification, achieving high accuracy in the optimization process and low level of error rate. A hybrid operation of Multi Support Vector Machine (SVM) with the Convolutional Neural Networks (CNN) have been carried out to attain effective results. The characteristics of the pre-processed face image are efficiently extracted by using LBP,PCAand Gray Level Occurrence Matrix (GLCM). For a trainable feature extractor CNN is used, and for a recognizer, Multi-SVM has been performed. The proposed system's performance is analysed with various human faces using the MATLAB tool. The results proved that the proposed method surpasses the earlier methods regarding high accuracy with low computation time and low error rate.

## REFERENCES

[1] C. Darwin, The Expression of the Emotions in Man and Animals. D. Appleton and Company, 1899.

[2] D. Keltner and P. Ekman, Handbook of Emotions, ch. 15 - Facial Expression of Emotion, pp. 151–249. Guilford Publications, Inc., 2nd ed., 2000.

[3] E. M. Provost, Y. Shangguan, and C. Busso, "Umeme: University of Michigan emotional McGurk effect data set," IEEE Transactions on Affective Computing, vol. 6, pp. 395–409, Oct 2015.

[4] R. E. Jack, O. G. Garrod, and P. G. Schyns, "Dynamic facial expressions of emotion transmit an evolving hierarchy of signals over time," Current Biology, vol. 24, no. 2, pp. 187 – 192, 2014.

[5] Lopes, A. T., de Aguiar, E., De Souza, A. F., & Oliveira-Santos, T. (2017). Facial expression recognition with convolutional neural networks: coping with few data and the training sample order. *Pattern Recognition*, *61*, 610-628.

[6] R. E. Jack, O. G. Garrod, and P. G. Schyns, "Dynamic facial expressions of emotion transmit an evolving hierarchy of signals over time," Current Biology, vol. 24, no. 2, pp. 187 – 192, 2014.

[7] M. W. Schurgin, J. Nelson, S. Iida, H. Ohira, J. Y. Chiao, and S. L. Franconeri, "Eye movements during emotion recognition in faces," Journal of Visualization, vol. 14, no. 13, 2014.

[8] I. Goodfellow, Y. Bengio, and A. Courville, "Deep learning." Book in preparation for MIT Press, 2016.

[9] Pilla Jr, V., Zanellato, A., Bortolini, C., Gamba, H. R., Borba, G. B., & Medeiros, H. Facial Expression Classification Using Convolutional Neural Network and Support Vector Machine.

[10] Azcarate, A., Hageloh, F., Van de Sande, K., & Valenti, R. (2005). Automatic facial emotion recognition. *Universiteit van Amsterdam*, 1-6.

[11] Cohen, I., Garg, A., & Huang, T. S. (2000, November). Emotion recognition from facial expressions using multilevel HMM. In *Neural information processing systems* (Vol. 2).

[12] Bartlett, M. S., Littlewort, G., Frank, M., Lainscsek, C., Fasel, I., & Movellan, J. (2006, April). Fully automatic facial action recognition in spontaneous behavior. In *Automatic Face and Gesture Recognition, 2006. FGR 2006. 7th International Conference on* (pp. 223-230). IEEE.

[13] Shan, C., Gong, S., & McOwan, P. W. (2009). Facial expression recognition based on local binary patterns: A comprehensive study. *Image and Vision Computing*, *27*(6), 803-816.

[14] Walecki, R.; Rudovic, O. Deep structured learning for facial expression intensity estimation. Image Vis. Comput. 2017, 259, 143–154.

[15] Al-Shabi, M., Cheah, W. P., & Connie, T. (2016). Facial Expression Recognition Using a Hybrid CNN-SIFT Aggregator. *arXiv preprint arXiv:1608.02833*.

[16] Hasani, B., & Mahoor, M. H. (2017, May). Spatio-temporal facial expression recognition using convolutional neural networks and conditional random fields. In *Automatic Face & Gesture Recognition (FG 2017), 2017 12th IEEE International Conference on* (pp. 790-795). IEEE.

[17] Yang, B., & Chen, S. (2013). A comparative study on the local binary pattern (LBP) based face recognition: LBP histogram versus LBP image. *Neurocomputing*, *120*, 365-379.

[18] Chamasemani, F. F., & Singh, Y. P. (2011, September). Multi-class support vector machine (SVM) classifiers--an application in hypothyroid detection and classification. In *Bio-Inspired Computing: Theories and Applications (BIC-TA), 2011 Sixth International Conference on* (pp. 351-356). IEEE.