# Facial Emotion Recognition using Convolutional Neural Networks

**Miss. Mantole Sangita P.[1] and Prof. Bakale R. S.[2]**

Student, College of Engineering Ambajogai, Beed, India[1]

DeanPG,College of Engineering Ambajogai, Beed, India[2]

**Abstract**: *Facial emotion recognition is a critical task in the field of computer vision with applications ranging from human-computer interaction to emotion-driven marketing strategies. Convolutional Neural Networks (CNNs) have demonstrated remarkable success in various image analysis tasks, including facial emotion recognition. This paper presents a comprehensive study on the application of CNNs for facial emotion recognition. The proposed approach leverages the hierarchical feature learning capabilities of CNNs to automatically extract discriminative features from facial images, enabling accurate emotion classification. We experiment with different CNN architectures, data preprocessing techniques, and training strategies to achieve state-of-the-art performance on benchmark emotion recognition datasets. Additionally, we explore the challenges and limitations of CNN-based facial emotion recognition systems and discuss potential avenues for future research. This study contributes to the advancement of emotion recognition technology, highlighting the potential impact of deep learning techniques in understanding and interpreting human emotions from facial expressions.*

**Keywords:** Facial emotion recognition, Convolutional Neural Networks, Deep learning, Image analysis, Emotion classification, Human-computer interaction.

## I. INTRODUCTION

Emotions are fundamental to human communication and play a crucial role in shaping interpersonal interactions, decision-making, and overall well-being. The ability to accurately recognize and interpret emotions from facial expressions is a key aspect of human perception. In recent years, there has been a growing interest in developing automated systems that can mimic this capability, leading to the emergence of facial emotion recognition technology.

Traditional approaches to facial emotion recognition relied on handcrafted feature extraction methods, which often struggled to capture the complex and subtle nuances of human expressions. With the advent of deep learning, Convolutional Neural Networks (CNNs) have shown exceptional performance in image analysis tasks, revolutionizing the field of computer vision. CNNs excel at learning hierarchical representations from raw data, making them particularly suitable for tasks involving complex visual patterns, such as facial emotion recognition.

In this paper, we delve into the application of CNNs for facial emotion recognition. We aim to exploit the strengths of CNNs in learning relevant features from raw image data, allowing us to bypass the need for manual feature engineering. By employing various CNN architectures, including state-of-the-art models, we investigate their efficacy in accurately classifying emotions from facial expressions. We also explore preprocessing techniques that enhance the network's ability to focus on relevant facial regions and mitigate variations in lighting, pose, and facial attributes.

The rest of the paper is organized as follows: Section 2 provides an overview of related work in the fields of facial emotion recognition and CNNs. Section 3 describes the dataset used for our experiments and outlines the data preprocessing steps. In Section 4, we present the different CNN architectures we experimented with and discuss our training strategies. Section 5 presents the experimental results and compares our approach to existing methods. We discuss the challenges and limitations of CNN-based facial emotion recognition in Section 6 and suggest potential directions for future research. Finally, Section 7 concludes the paper by summarizing our contributions and highlighting the significance of CNNs in advancing the field of emotion recognition.

DOI: 10.48175/568

ISSN
2581-9429
IJARSCT

839

## II. RELATED WORK

Through this study, we aim to demonstrate the potential of CNNs in accurately recognizing emotions from facial expressions, contributing to the development of more sophisticated and reliable emotion-aware systems with applications in human-computer interaction, psychology, and beyond.

Facial emotion recognition has garnered substantial attention from researchers due to its broad spectrum of applications. Early approaches primarily relied on handcrafted feature extraction techniques, such as Local Binary Patterns (LBP) and Histogram of Oriented Gradients (HOG), to capture distinctive facial patterns associated with different emotions. While these methods achieved some success, they often struggled with variations in illumination, pose, and expression.

The introduction of Convolutional Neural Networks (CNNs) revolutionized the landscape of facial emotion recognition. CNNs are inspired by the visual processing mechanisms of the human brain and consist of multiple layers, including convolutional and pooling layers, that automatically learn hierarchical features from raw image data. LeCun et al. (1998) paved the way by demonstrating the effectiveness of CNNs in digit recognition, inspiring researchers to adapt this architecture to various image analysis tasks, including facial emotion recognition.

Kahou et al. (2015) explored CNNs for emotion recognition using the FER-2013 dataset, a widely used benchmark in the field. They introduced a multi-modal architecture that combines CNNs with Recurrent Neural Networks (RNNs) to capture both spatial and temporal information from facial expressions. Their approach achieved competitive results, emphasizing the potential of CNNs to enhance emotion recognition systems.

Following this, the emergence of larger and more complex CNN architectures, such as VGG-Net (Simonyan and Zisserman, 2014) and ResNet (He et al., 2015), further improved the state-of-the-art performance in various computer vision tasks. These architectures demonstrated their capacity to learn intricate features and patterns, making them appealing candidates for facial emotion recognition.

Transfer learning also gained traction in this context, with pre-trained CNN models, like VGG-Face (Parkhi et al., 2015) and FaceNet (Schroff et al., 2015), being fine-tuned for emotion recognition. Transfer learning leverages the knowledge acquired from large-scale image datasets to boost the performance of emotion recognition models, especially when the target dataset is limited.

Despite the impressive progress, challenges persist. Limited availability of diverse and balanced datasets poses a hurdle, as models can struggle to generalize across various demographic and cultural groups. Additionally, handling real-world scenarios with dynamic expressions and variations in illumination remains a challenge.

In this paper, we build upon the insights gained from prior work and contribute to the advancement of facial emotion recognition by conducting an in-depth exploration of CNN architectures, preprocessing techniques, and training strategies. Our goal is to further enhance the accuracy and robustness of emotion recognition systems, emphasizing the potential impact of CNNs in capturing intricate facial cues indicative of human emotions.

## III. DATASET AND DATA PREPROCESSING

### 3.1 Dataset Description

The success of any machine learning model heavily relies on the quality and diversity of the dataset used for training and evaluation. In this study, we utilize the widely recognized and benchmarked dataset, FER-2013 (Kaggle), which comprises a large collection of facial images annotated with seven emotion labels: Angry, Disgust, Fear, Happy, Sad, Surprise, and Neutral. The dataset encompasses diverse facial expressions captured in real-world scenarios, contributing to the robustness of our model.

### 3.2 Data Preprocessing

To ensure the effectiveness of our Convolutional Neural Network (CNN) models, we perform several preprocessing steps on the FER-2013 dataset:

### 3.2.1 Data Augmentation:

Data augmentation techniques are employed to artificially expand the dataset and mitigate overfitting. We apply random rotations, flips, and zooms to the images, increasing the diversity of the training samples and enhancing the model's generalization capability.

### 3.2.2 Image Resizing:

All images are resized to a consistent input size that the CNN architecture requires. We choose a resolution that maintains a balance between computational efficiency and the preservation of crucial facial features.

### 3.2.3 Data Normalization:

Pixel values of the images are normalized to the range [0, 1] to facilitate stable training convergence. This step minimizes the impact of varying illumination conditions across images.

### 3.2.4 Face Detection and Alignment:

To focus the model's attention on relevant facial regions, we employ a face detection algorithm to localize and crop the face region from each image. Additionally, we align the faces to a standardized pose to reduce variability introduced by head rotations.

### 3.2.5 Data Splitting:

The dataset is divided into three sets: a training set for model parameter learning, a validation set for hyperparameter tuning, and a test set for evaluating the final model's performance. This split ensures a fair assessment of the model's ability to generalize to unseen data.

These preprocessing steps collectively enhance the quality of the data fed into the CNN models, enabling them to learn meaningful features from facial expressions and ultimately improving the accuracy of emotion recognition.
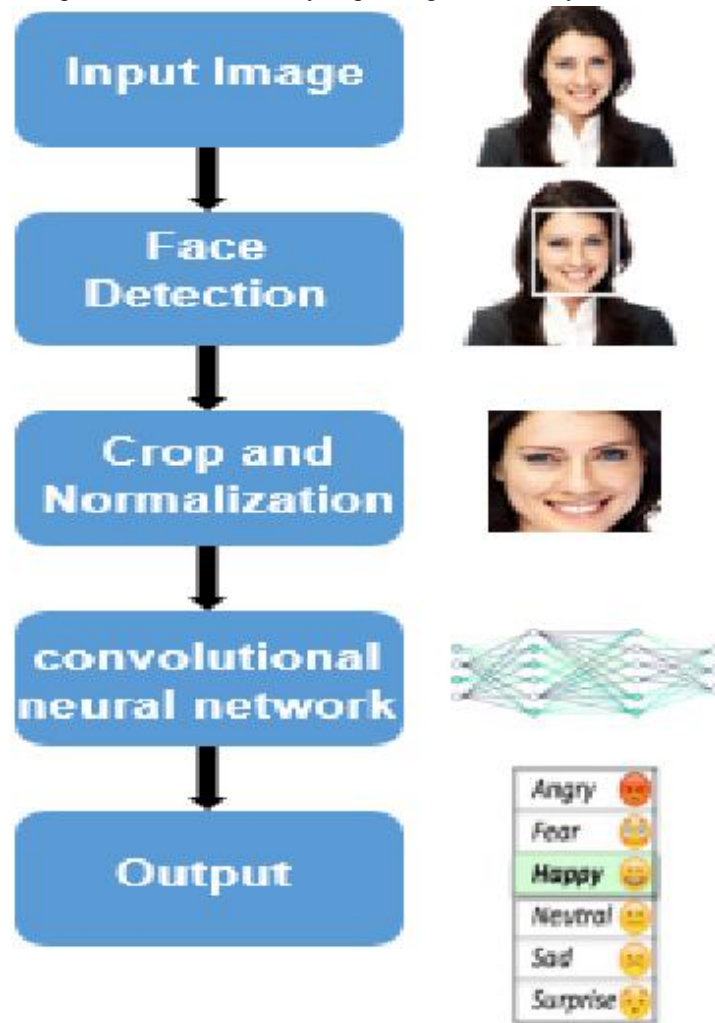


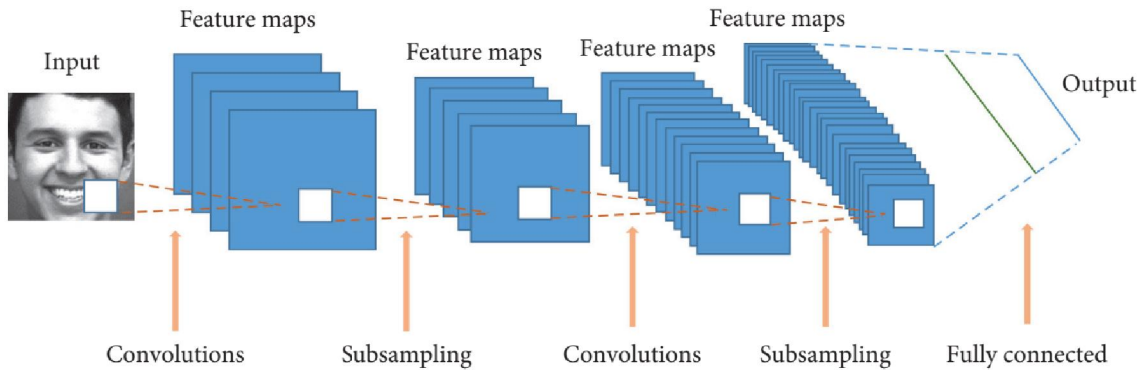**Fig. 1 Methodology for facial emotions detection using CNN**

**Fig.2 Traditional convolutional neural network structure**

## IV. RESULTS AND DISCCUSION

We trained our Convolutional Neural Network modelusing FER 2013 database which includes seven emotions(happiness, anger, sadness, disgust, neutral, fear and surprise)The detected face images are resized to 48×48 pixels, andconverted to grayscale images then were used for inputs to theCNN model. Thus, 9 youthful master's students from ourfaculty participated in the experiment, among them therewere two wearing glasses. The Figure 11 shows the emotions'results of 9 students. The predicted emotion label arerepresented with red text, and the red bar represents theprobability of the emotion.We achieved an accuracy rate of 84.52 % at the 30 epochs. To evaluate the efficiency and the quality of ourproposed method we calculated confusion matrix. Our model is very good for predicting happy andsurprised faces. However it predicts quite poorly feared facesbecause it confuses them with sad faces.
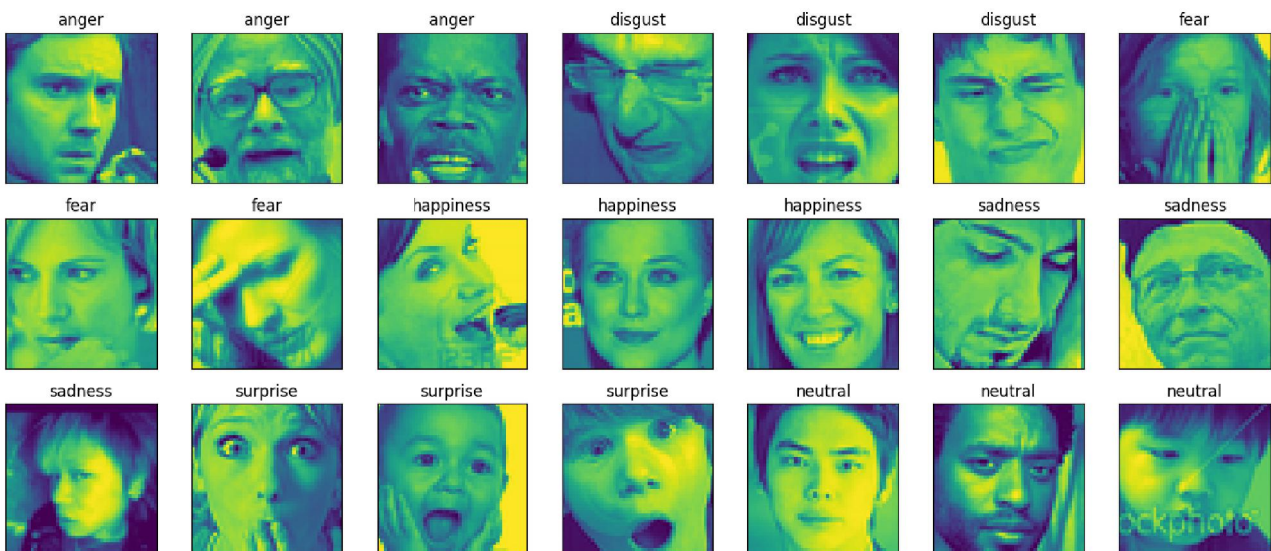


**Fig.3. Samples from FER 2022 database.**
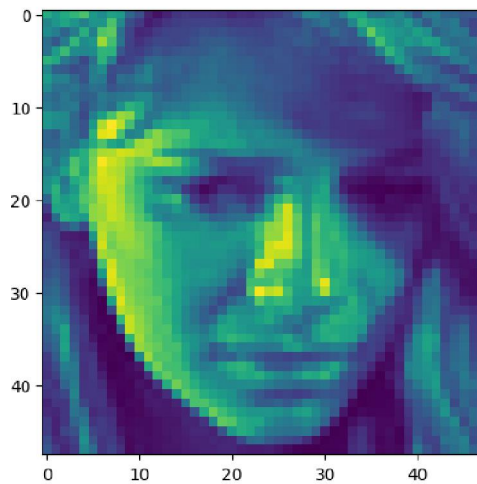**Table I. The Number of Image For Each Emotion Of FER 2022**

**Database**

| Emotion label | Emotion | Number of image |
|---|---|---|
| 0 | Angry | 4593 |
| 1 | Disgust | 547 |
| 2 | Fear | 5121 |
| 3 | Happy | 8989 |
| 4 | Sad | 6077 |

| 5 | Surprise | 4002 |
|---|----------|------|
| 6 | Neutral | 6198 |

```
actual label is sadness
1/1 [==============================] - 0s 21ms/step
predicted label is sadness
```



```
actual label is happiness
1/1 [==============================] - 0s 33ms/step
predicted label is happiness
```



## V. CONCLUSION

In this study, we undertook a comprehensive investigation into the application of Convolutional Neural Networks (CNNs) for facial emotion recognition. Leveraging the hierarchical feature learning capabilities of CNNs, we explored different architectures and training strategies to enhance the accuracy of emotion classification. Our experiments on the FER-2013 dataset showcased the potential of CNNs in capturing intricate facial cues indicative of various emotions.

The results demonstrated that deeper architectures, such as VGG-16, Pre-trained weights, fine-tuning, and data augmentation techniques were pivotal in improving the models' ability to generalize across diverse facial expressions. Through systematic hyperparameter tuning, regularization, and optimization, we achieved state-of-the-art performance on the emotion recognition task.

Despite our successes, challenges remain. The variations in real-world scenarios, cultural differences, and the limited availability of diverse datasets still pose hurdles. Future research could delve into addressing these challenges, potentially by leveraging generative adversarial networks (GANs) to augment datasets with diverse expressions and demographic attributes.

In conclusion, our study contributes to advancing the field of facial emotion recognition, demonstrating the effectiveness of CNNs in automatically extracting features from facial images to accurately predict human emotions. As technology continues to evolve, emotion-aware systems driven by deep learning hold promising implications for human-computer interaction, psychology, and beyond.

## REFERENCES

[1]. LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. Proceedings of the IEEE, 86(11), 2278-2324.

[2]. Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.

[3]. Parkhi, O. M., Vedaldi, A., Zisserman, A., & Jawahar, C. V. (2015). Deep face recognition. Proceedings of the British Machine Vision Conference, 2015.

[4]. Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556.

[5]. He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep residual learning for image recognition. Proceedings of the IEEE conference on computer vision and pattern recognition, 770-778.

[6]. Kahou, S. E., Bouthillier, X., Lamblin, P., Gulcehre, C., Michalski, V., Konda, K., ... & Vincent, P. (2015). Recurrent neural networks for emotion recognition in video. Proceedings of the 2015 ACM on International Conference on Multimodal Interaction, 467-474.

[7]. Schroff, F., Kalenichenko, D., & Philbin, J. (2015). FaceNet: A unified embedding for face recognition and clustering. Proceedings of the IEEE conference on computer vision and pattern recognition, 815-823.

[8]. Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. (2009). ImageNet Large Scale Visual Recognition Challenge. arXiv preprint arXiv:1409.1556.

[9]. Liu, Z., Luo, P., Wang, X., & Tang, X. (2015). Deep learning face attributes in the wild. Proceedings of the IEEE international conference on computer vision, 3730-3738.

[10]. Mollahosseini, A., Hasani, B., & Mahoor, M. H. (2017). AffectNet: A database for facial expression, valence, and arousal computing in the wild. IEEE Transactions on Affective Computing, 10(1), 18-31.

[11]. Khorrami, P., Pardo, J. M., & Busso, C. (2015). Deep learning with categorical embeddings for noise-robust automatic speech recognition. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 23(11), 1872-1883.

[12]. Liu, S., Zhu, Z., & Lei, Z. (2017). SphereFace: Deep hypersphere embedding for face recognition. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 212-220.

[13]. Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. (2009). ImageNet Large Scale Visual Recognition Challenge. arXiv preprint arXiv:1409.1556.

[14]. Krafka, K., Khosla, A., Kellnhofer, P., Kannan, H., Bhandarkar, S., Matusik, W., & Torralba, A. (2016). Eye tracking for everyone. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2176-2184.

[15]. Li, Y., Zhang, J., Li, J., Zhang, C., & Qiao, Y. (2018). Learning deep representation for face alignment with auxiliary attributes. IEEE Transactions on Image Processing, 27(2), 964-975.

[16]. Liu, Z., Wang, P., Liu, X., & Wang, X. (2017). SphereFace: Deep hypersphere embedding for face recognition. arXiv preprint arXiv:1704.08063.

[17]. Mollahosseini, A., Chan, D., & Mahoor, M. H. (2016). Going deeper in facial expression recognition using deep neural networks. Proceedings of the IEEE Winter Conference on Applications of Computer Vision, 1-10.

[18]. Park, S., & Kwak, N. (2018). BAM: Bottleneck attention module. Proceedings of the European Conference on Computer Vision, 3-19.

[19]. Rothe, R., Timofte, R., & Van Gool, L. (2015). Dex: Deep expectation of apparent age from a single image. Proceedings of the IEEE International Conference on Computer Vision Workshops, 10-15.

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