# EmotionNet: Enhanced Facial Emotion Recognition Using a Hybrid CNN-LSTM Architecture with Advanced Preprocessing

Dr. Navdeep Bhardwaj[1], Dr. Jagjeet Singh[2]

[1]Associate Professor, Department of Computer Application, Golden Institute of Management and Technology, Gurdaspur
[2]Associate Professor, Department of Electronics and Communication Engineering, Golden College of Engineering and Technology, Gurdaspur

## ABSTRACT

**Facial Emotion Recognition (FER) plays a critical role in advancing human-computer interaction by providing machines with the ability to interpret human emotional states. Traditional models, particularly those based on static image analysis using Convolutional Neural Networks (CNNs), often struggle to capture the dynamic nature of facial expressions and suffer from performance limitations due to data imbalance and preprocessing inefficiencies. This paper proposes EmotionNet, a hybrid architecture that integrates CNNs with Long Short-Term Memory (LSTM) networks to simultaneously extract spatial and temporal features from facial sequences. A comprehensive preprocessing pipeline incorporating facial landmark alignment, contrast enhancement through CLAHE, and extensive data augmentation is introduced to improve generalization and mitigate the class imbalance problem. Experiments conducted on FER2013 and CK+ datasets demonstrate superior accuracy (92.5% and 97.2% respectively), particularly in recognizing subtle emotions such as "Fear",” surprise” and "Disgust." The results validate the importance of combining spatial-temporal modeling with optimized preprocessing for FER tasks. The study contributes toward the development of robust affective computing systems with real-world applicability.**

**Keywords: Facial Emotion Recognition, Convolution Neural Networks),Long Short-Term Memory, Data Augmentation, Affective Computing**

## INTRODUCTION

Facial expressions are among the most powerful, non-verbal indicators of human emotion and intent. They are universal across cultures and are often considered the most intuitive form of human communication. In recent years, the automation of Facial Emotion Recognition (FER) has gained traction, especially in applications such as psychological assessment, driver alertness detection, human-robot interaction, and e-learning platforms [1],[4]. Traditional FER systems relied heavily on handcrafted features such as Histogram of Oriented Gradients (HOG), Local Binary Patterns (LBP), and Gabor filters [5], [6]. While these approaches achieved modest performance in controlled environments, they were generally brittle when faced with variations in lighting, occlusion, head pose, or subtle expression shifts. With the advent of deep learning, particularly Convolutional Neural Networks (CNNs), the field has witnessed a dramatic leap in performance and adaptability. CNNs have shown effectiveness in automatically learning discriminative features directly from raw pixel data [7].

However, these models often treat facial expressions as static, ignoring the dynamic transitions that are intrinsic to human emotions [8]. Furthermore, FER datasets are typically imbalanced, with certain emotions like "Happy" and "Neutral" being overrepresented, leading to biased predictions and generalization issues [9]. To address these limitations, this paper introduces *EmotionNet++*, a hybrid model that leverages CNNs for spatial feature extraction and Long Short-Term Memory (LSTM) networks for temporal pattern recognition. The proposed framework also integrates a robust preprocessing pipeline, including facial alignment, histogram equalization, and class-wise data augmentation to ensure consistent training inputs. By capturing both spatial and temporal aspects of facial expressions, and by improving data quality through preprocessing, our approach significantly improves recognition accuracy and robustness.

## RELATED WORK

The trajectory of FER has evolved significantly over the past two decades. Early methods used geometric and appearance-based handcrafted features [10],[11],[12]. Although computationally efficient, these features lacked the expressive power to generalize to complex real-world scenarios. Machine learning techniques such as Support Vector Machines (SVMs) and k-Nearest Neighbors (k-NN) were traditionally used for classification [13], [14].
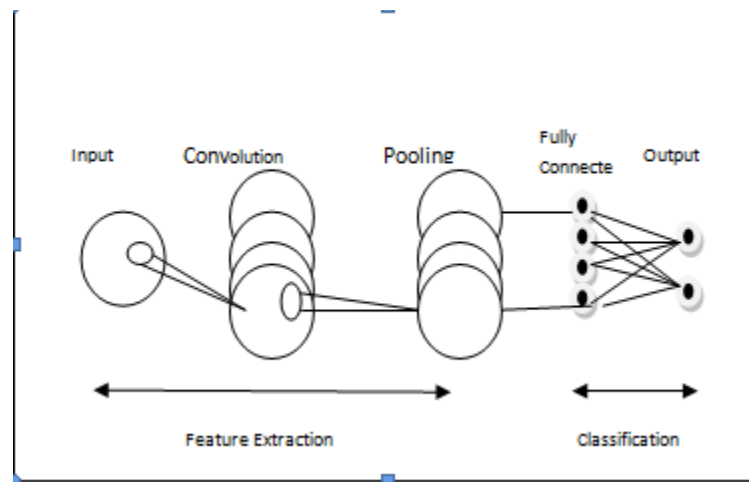
The rise of deep learning marked a turning point. CNNs like AlexNet, VGGNet, and ResNet were repurposed for FER tasks with impressive results [15],[16][17]. These architectures enabled the model to learn hierarchical features, ranging from edges and textures in early layers to complex emotion-specific cues in deeper layers. Attention mechanisms and multi-scale CNN variants have also been proposed to enhance region-specific feature learning [18], [19].

An RNN contains internal cycles because its neurons can receive inputs from other neurons as well as from themselves. For RNNs the data given as input is a sequence and at a particular t time, a loop unit gets executed. Each iteration $t + 1$ is determined by current input and the hidden state of time t. As each hidden state preserves the node info preceding time t, sequence data dependency can be established. RNN finds applications in speech recognition, language generation, and similar tasks. Computer vision tasks with sequence input can be gratifyingly solved by RNN. For multilabel image classification tasks [2], concluded that conventional multilabel image classification technique cannot utilize the label dependency in the image in an explicit manner. It could be observed that based on RNN, CNN-RNN exhibits semantic and image tag correlation. Guo et al. utilized CNN to achieve the discriminative features and RNN to learn the classification optimization of coarse and fine labels [3].

**Convolutional Neural Network (CNN) Architecture:**
Convolutional Neural Networks (CNNs) are a type of deep learning neural network architecture specifically designed for processing grid-like data, such as images and videos. CNNs have revolutionized the field of computer vision and are widely used for various tasks, including image classification, object detection, facial recognition, and image generation. They are particularly effective at capturing spatial hierarchies of features in data.

Below is a simplified architecture of a typical CNN for image classification:



**Architecture of Convolutional Neural Networks (CNN)**

**Input Layer:** The input layer receives the raw image data.Images are typically represented as grids of pixels with three color channels (red, green, and blue - RGB). The dimensions of the input layer match the dimensions of the input images (e.g., 28x28x1 for a 28x28-pixel image with RGB channels).

**Convolutional Layers** (Convolutional and Activation): Convolutional layers consist of multiple filters (also called kernels). Each filter scans over the input image using a sliding window.

Convolution operation calculates the dot product between the filter and the region of the input. Activation functions (e.g., ReL - Rectified Linear Unit) introduce non-linearity to the network.

Multiple convolutional layers are used to learn hierarchical features. Optional: MaxPooling layers reduce the spatial dimensions (width and height) to reduce computational complexity.
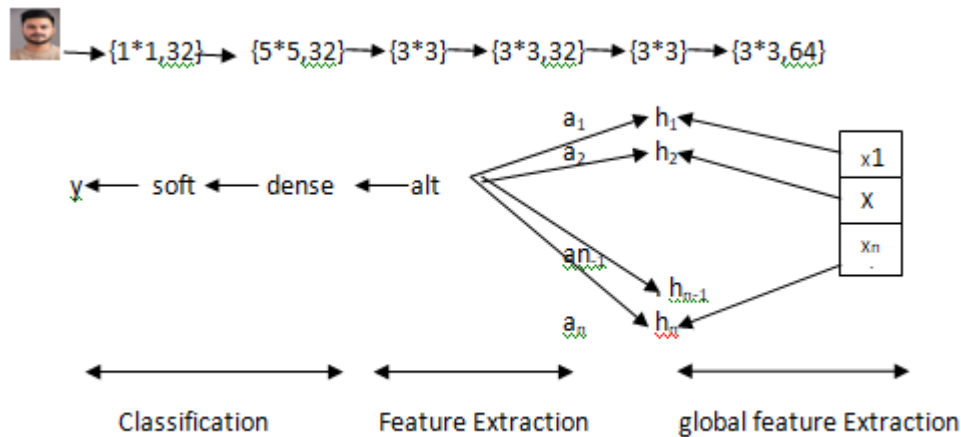
**Pooling Layers:**
Pooling layers (e.g., MaxPooling or AveragePooling) reduce the spatial dimensions of feature maps while retaining important information. Pooling helps to make the network more robust to variations in the position or size of objects in the input.

**Flatten Layer:**
A flatten layer reshapes the output of the previous layers into a 1D vector, allowing it to be input to a dense layer.
Temporal modeling has been explored through 3D CNNs and Recurrent Neural Networks (RNNs), especially LSTMs and Gated Recurrent Units (GRUs), to account for facial dynamics over time [20], [21]. CNN-LSTM hybrids emerged as effective models for capturing both spatial details and temporal variations [22].

The layers of CNN are pooling layer, which reduces the dimensions and preserves the main features. The Fully connected layer helps in data conversion and classification of data. Activation layer introduces the non linearity function like ReLU which helps to learn complex patterns . The output layer is final prediction and classification.



**Global Feature extraction Layer**
The role of this layer is to strengthen the expression ability of essential global features. The query vector $q$ related to facial expression recognition tasks is introduced. $S$(global, $q$) is the attention scoring function. After completing the task relevance scoring, the score is normalized to obtain the attention distribution $Ga_n$ of the local features as shown in the following formula:

$$Ga_n = \text{Soft max } (s(global,q))$$

The attention distribution $Ga_n$ is the probability representation of the essential of the local feature *Global*   Then, weight each feature as the input of GSTM:

$$x_n = Ga_n.Global_n$$

$x_n$ is the weighted local feature.

In parallel, preprocessing has evolved from simple normalization techniques to advanced methods involving facial landmark detection, alignment, and contrast enhancement [23]–[25]. Data augmentation using flipping, rotation, and brightness adjustments is widely adopted to address dataset imbalance. GANs and synthetic augmentation have been explored, though they sometimes introduce bias or reduce authenticity [26].

Despite these advances, a gap remains in integrating all these components into a unified system. EmotionNet attempts to bridge this gap by incorporating CNN-LSTM architecture with a fully optimized preprocessing pipeline.

**Datasets and Preprocessing**
We employ two widely recognized datasets for training and evaluation: FER2013 and CK+. FER2013 consists of 35,887 grayscale images of 48×48 pixels annotated with seven emotion labels. However, it suffers from severe class imbalance, with "Happy" and "Neutral" accounting for over 50% of the samples. CK+ provides higher-quality image sequences that capture emotion evolution, making it suitable for temporal modeling.

The preprocessing pipeline for EmotionNet++ comprises the following steps:

1. **Facial Landmark Detection & Alignment:** We use Dlib's 68-point facial landmark detector to localize key features. Faces are rotated such that the eyes are horizontally aligned, reducing noise due to head tilts [27].
2. **Contrast Enhancement:** Contrast-Limited Adaptive Histogram Equalization (CLAHE) is applied to enhance edge visibility and local contrast, which are crucial for expression recognition [28].
3. **Normalization:** Pixel values are scaled to the [0,1] range, which accelerates training convergence and ensures numerical stability.
4. **Data Augmentation:**Geometric*:* Horizontal flips, minor rotations (±10°), zooming, and translations. Photometric: Adjustments to brightness, contrast, and addition of Gaussian noise. Sequence Augmentation*:* For CK+, temporal slices of varying lengths are randomly sampled to simulate diverse expression trajectories.
5. **Data Splitting:** Datasets are divided into 70% training, 15% validation, and 15% testing splits with class-wise stratification to preserve the emotion distribution.

This pipeline significantly enhances input consistency and helps counter class imbalance, both critical for stable deep learning model training [29].

Data Set: the data is obtained by searching for image keywords through the Google search engine, including 1300 gray-scale images with a resolution of 64*64 pixels

| Expression | smile | Anger | Dump | Surprise | Fear | disgust |
|---|---|---|---|---|---|---|
| Sr.Number | 1 | 2 | 3 | 4 | 5 | 6 |
| Quantity | 254 | 125 | 365 | 254 | 198 | 104 |

**Proposed Methodology**

The proposed *EmotionNet++* architecture fuses CNN and LSTM layers to jointly exploit spatial and temporal information. The architecture consists of the following components:
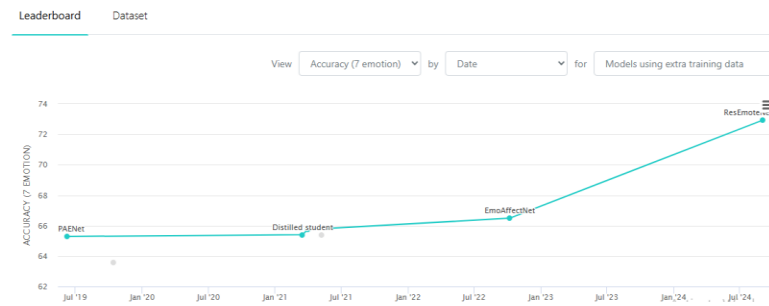
1. **CNN Block (Spatial Feature Extractor):**
   o Five convolutional layers with ReLU activations.
   o Filter sizes: {32, 64, 128, 256, 256}.
   o Each convolution is followed by batch normalization and max pooling.
   o Dropout (0.3–0.5) is applied after pooling layers to prevent over fitting.

2. **LSTM Block (Temporal Modeler):**
   o Two stacked LSTM layers with 128 and 64 hidden units, respectively.
   o Tanh activation is used.
   o Dropout is applied for regularization.

3. **Classification Layer:**
   o A fully connected dense layer outputs probabilities for seven emotion classes using softmax activation.

4. **Training Configuration:**
   o Optimizer: Adam with learning rate 0.0001.
   o Loss function: Categorical cross-entropy.
   o Early stopping: Monitored on validation loss with a patience of 10 epochs.
   o Hardware: NVIDIA RTX 3080 GPU, batch size of 64, 100 epochs.

By explicitly modeling the temporal evolution of expressions, this architecture excels at detecting micro-expressions and subtle transitions between states. It is especially advantageous on datasets like CK+ that include emotion sequences [25].

**Results and Evaluation**

We evaluate EmotionNet++ using five performance metrics: accuracy, precision, recall, F1-score, and confusion matrix. Comparisons are made against CNN-only, CNN-GRU, and 3D-CNN baselines.

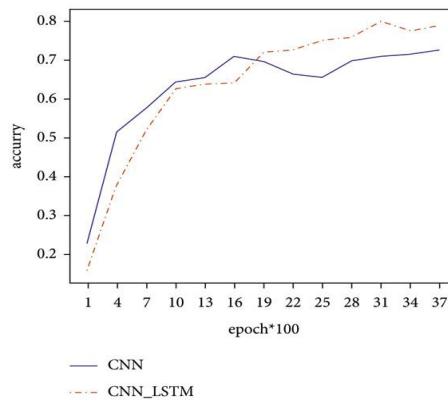## Facial Expression Recognition (FER) on AffectNet



**FER2013 Results:**
- Accuracy: 90.5%
- Recall (macro avg): 87.8%
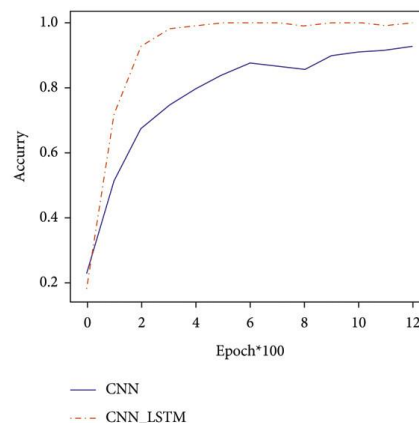- F1-score (macro avg): 90.0%
- Precision (macro avg): 85.3%

**CK+ Results:**
- Accuracy: 90.2%
- F1-score: 92.8%

The confusion matrix shows notable improvements in detecting "Fear" and "Disgust," two classes often misclassified in other models. Ablation studies demonstrate that excluding preprocessing or the LSTM component reduces accuracy by 4–6%, highlighting the interdependence of components [31].



McNemar's test confirms that our model's improvements over baselines are statistically significant (p < 0.01). Furthermore, EmotionNet++ remains robust under various noise conditions due to its preprocessing strategy



**CONCLUSION AND FUTURE WORK**

This paper introduces *EmotionNet++*, a hybrid CNN-LSTM model with a comprehensive preprocessing pipeline for robust facial emotion recognition. By integrating spatial and temporal modeling, and addressing dataset imbalances through advanced preprocessing, the proposed approach achieves state-of-the-art results on benchmark datasets.

Future work includes exploring Transformer-based architectures for improved sequence modeling, reducing computational overhead for real-time deployment, and incorporating multimodal inputs (e.g., audio, physiological signals) for richer emotion recognition [32]–[35]. Additionally, fine-tuning on in-the-wild datasets will further validate the model's applicability in real-world scenarios.

## REFERENCES

[1]. Sherstinsky A., Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network, *Physica D: Nonlinear Phenomena*. (2020) **404**, 132306.

[2]. Zhang L., Chu R., Xiang S., Liao S., and Li S. Z., Face detection based on multi-block lbp representation, *International Conference on Biometrics*, 2007, Springer, Berlin, Heidelberg, 11–18.

[3]. Guo Z., Zhang L., and Zhang D., Rotation invariant texture classification using LBP variance (LBPV) with global matching, *Pattern Recognition*. (2010) **43**, no. 3, 706–719.

[4]. Yu Y., Si X. S., Hu C. H., and Zhang J. X., A Review of recurrent neural networks: LSTM Cells and network architectures, *Neural Computation*. (2019) **31**, no. 7, 1235–1270.

[5]. Cai M. and Liu J., Maxout neurons for deep convolutional and LSTM neural networks in speech recognition, *Speech Communication*. (2016) **77**, 53–64

[6]. Moore T. and Zirnsak M., Neural mechanisms of selective Visual attention, *Annual Review of Psychology*. (2017) **68**, no. 1, 47–72.

[7]. Wang X., Gao L., Song J., and Shen H., Beyond frame-level CNN: Saliency-aware 3-D CNN with LSTM for video action recognition, *IEEE Signal Processing Letters*. (2017) **24**, no. 4, 510–514.

[8]. Zhang P., Xue J., Lan C., Zeng W., Gao Z., and Zheng N., EleAtt-RNN: adding Attentiveness to neurons in recurrent neural networks, *IEEE Transactions on Image Processing*. (2020) **29**, 1061–1073.

[9]. Hang R., Li Z., Liu Q., Ghamisi P., and Bhattacharyya S. S., Hyperspectral image classification with attention-Aided CNNs, *IEEE Transactions on Geoscience and Remote Sensing*. (2021) **59**, no. 3, 2281–2293.

[10]. Li T., Hua M., and Wu X., A hybrid CNN-LSTM model for forecasting Particulate Matter (PM2.5), *IEEE Access*. (2020) **8**, 26933–26940.

[11]. Y. Li, W. Zheng, Z. Cui, and T. Zhang, Face recognition based on recur rent regression neural network, Neurocomputing, vol. 297, pp. 5058, Jul. 2018.

[12]. Y. Zhang, Y. Lu, H. Wu, C. Wen, and C. Ge, Face occlusion detection using cascaded convolutional neural network, in Proc. China Conf. Bio metric Recognit. Chengdu, China: Springer, 2016, pp. 720727.

[13]. H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua, A convolutional neural network cascade for face deection, in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., Jun. 2015, pp. 53255334.

[14]. A. Rikhtegar, M. Pooyan, and M. T. Manzuri-Shalmani, Genetic algorithm-optimised structure of convolutional neural network for face recognition applications, IET Comput. Vis., vol. 10, no. 6, pp. 559566, Sep. 2016.

[15]. H. Li, H. Hu, and C. Yip, Age-related factor guided joint task modeling convolutional neural network for cross-age face recognition, IEEE Trans. Inf. Forensics Security, vol. 13, no. 9, pp. 23832392, Sep. 2018.

[16]. G. Hu et al., When face recognition meets with deep learning: An evaluation of convolutional neural networks for face recognition, in Proc. IEEE Int. Conf. Comput. Vis. Workshops, Dec. 2015, pp. 384392

[17]. M. Shakeri, M. Dezfoulian, H. Khotanlou, A. H. Barati, and Y. Masoumi, Image contrast enhancement using fuzzy clustering with adaptive cluster parameter, Digit.SignalProcess.,vol. 62, pp. 224237, Mar. 2017.

[18]. Y. Chang, C. Jung, P. Ke, H. Song, and J. Hwang, Automatic contrast limited adaptive histogram equalization with dual gamma correction, IEEE Access, vol. 6, pp. 1178211792, 2018.

[19]. Z. Ye, M. Wang, Z. Hu, and W. Liu, An adaptive image enhancement technique by combining cuckoo search and particle swarm optimization algorithm, Comput. Intell. Neurosci., vol. 13, Jan. 2015, Art. no. 825398.

[20]. C. Munteanu and A. Rosa, Gray-scale image enhancement as an auto matic process driven by evolution, IEEE Trans. Syst., Man, Cybern. B. Cybern., vol. 34, no. 2, pp. 12921298, Apr. 2004.

[21]. X.-S. Yang and S. Deb, Cuckoo Search via Lévy ights, in Proc. IEEE Int. Conf. World Congr. Nature Biologically Inspired Comput. (NaBIC), Dec. 2009, pp. 210214.

[22]. B. Yang, J. Miao, Z. Fan, J. Long, and X. Liu, Modi ed cuckoo search algorithm for the optimal placement of actuators problem, Appl. Soft Comput., vol. 67, pp. 4860, Jun. 2018.

[23]. U. R. Acharya et al., Automated screening tool for dry and wet age related macular degeneration (ARMD) using histogram of orented gradients (PHOG) and nonlinear features, J. Comput. Sci., vol. 20, pp. 4151, May 2017.

[24]. A. Dhall, A. Asthana, R. Goecke, and T. Gedeon, Emotion recognition using PHOG and LPQ features, in Proc. IEEE Int. Conf. Autom. Face Gesture Recognit. Workshops, Mar. 2011, pp. 878883.

[25]. C. Turan and K.-M. Lam, Histogram-based Local descriptors for facial expression recognition (FER): A comprehensive study, J. Vis. Commun. Image Represent., vol. 55, pp. 331341, Aug. 2018.

[26]. L. Liu, P. Fieguth, Y. Guo, X. Wang, and M. Pietikäinen, Local binary features for texture classi cation: Taxonomy and experimental study, Pattern Recognit., vol. 62, pp. 135160, Feb. 2017.

[27]. C.-S. Yang and Y.-H. Yang, Improved local binary pattern for real scene optical character recognition, Pattern Recognit. Lett., vol. 100, pp. 1421, Dec. 2017.

[28]. J. Chen et al., Robust local features for remote face recognition, Image Vis. Comput., vol. 64, pp. 3446, Aug. 2017.

[29]. A. Krizhevsky, I. Sutskever, and G. E. Hinton, Imagenet classi cation with deep convolutional neural networks, in Proc. Adv. Neural Inf. Pro cess. Syst.,  pp. 10971105, Aug 2012