# International Journal of Research Publication and Reviews

# Real-Time Facial Emotion Recognition

## *Rajashekar Potharaju[1], T. Malathi [2]*

*[1,] P.G. Research Scholar, Department of MCA-Data Science, Aurora Deemed To Be University, Hyderabad, India*
*[2]Assistant Professor, Department of MCA, Aurora Deemed To Be University, Hyderabad, Telangana, India*
Email: [1]rajashekarvarma702@gmail.com, [2]malathi@aurora.edu.in

## ABSTRACT

Human-Computer Interaction (HCI) development is increasingly based on the capacity of systems to perceive and react to human emotional state. This paper presents the design and deployment of a robust, real-time facial emotion detection system and aims at closing the human affective expression-computation understanding gap. Our primary contribution is a dense structure that not only accurately identifies human emotions but also expresses this explanation in synthesized speech and offers a more natural and interactive user interface.

To achieve this, we built and trained a deep Convolutional Neural Network (CNN) on a varied database of facial expressions such that the model was able to identify seven universal emotions accurately: anger, disgust, fear, happiness, sadness, surprise, and neutrality. The system employs OpenCV for real-time capture and processing of the video stream and Haar Cascade classifiers for face detection effectively. The found facial region of interest is then processed by our trained CNN. The second innovation of this work is in using a text-to-speech engine, which gives real-time, vocalized feedback of the classified emotion. This functionality makes the system shift from a passive viewer to an active agent of interaction. Performance testing shows the high classification accuracy of the model and the efficacy of the system in real time. This book is a seminal contribution to affective computing, with profound implications for assistive technology, interactive learning spaces, and more effective user experience design, through the development of systems that are not only smart but emotionally smart as well.

*Keywords: Emotion Recognition, Deep Learning, Computer Vision, Convolutional Neural Network (CNN), Real-Time Systems, Affective Computing, Human-Computer Interaction (HCI)*

## Introduction

With the rapidly evolving world of computers, the border line between human and machine interaction is quickly being outdated. The core of modern Human-Computer Interaction (HCI) has shifted from technical command-based systems to developing intuitive, adaptive, and user-state-sensitive systems. The paradigm shift has its focus on the field of Affective Computing, an inter-disciplinary research field that is concerned with developing systems and machines that are capable of perceiving, understanding, processing, and emulating human affects. As most human communication consists of non-verbal information, being able to perceive and understand emotions forms the foundation of the development of truly smart and empathic computational systems.

The human face is one of the most immediate and compelling ways of conveying emotional states. The face conveys high-bandwidth information rich in emotion and intentions. Facial Emotion Recognition (FER) is thus becoming a critical research space with important implications for a variety of applications, from learning systems and mental health monitoring to driver safety systems and socially assistive robots. The advent of deep learning, or Convolutional Neural Networks (CNNs), has transformed computer vision to enable high-accuracy and high-robustness models for Facial Emotion Recognition (FER) that can surpass traditional machine learning techniques.

But although much progress has been achieved in the reliable detection of facial emotions, most systems currently available are passive spectators. They might take emotional information or present a label on screen but generally neglect to round out the cycle of interaction by presenting their comprehension in a readable form to the user. This one-way flow of information makes a transactional rather than an interactional experience, shutting the door to richer and more profound interaction. The failure of return feedback is an essential lack between an internal "awareness" in a system and its external potential for sage interaction.

This work closes this gap by describing the design, prototyping, and testing of an actual facial emotion recognition system that not only identifies a user's emotional state but provides immediate vocal feedback. Our system accomplishes this by using a rigorously trained CNN paired with a text-to-speech (TTS) engine that takes the passive emotion recognition and turns it into an active two-way interaction. The long-term vision of this work is to create a more natural and conversational HCI in which the system is able to capture the emotional state of the user and thus make the user feel understood. This paper discusses our overall system architecture, from collecting data to training models to finally executing in real-time. We demonstrate that this closed-

loop approach, with high-fidelity vision augmented with audio feedback, is a pivotal step in the design of more emotionally responsive and sensitive computational companions.

The remainder of this paper is structured as follows: Section 2 provides the methodology, including the data, our Convolutional Neural Network design, and the system implementation. Section 3 presents the experimental findings, including the performance metrics of the model and classification accuracies. Section 4 is results discussion and implications, as well as the system limitations. Finally, Section 5 finishes the paper with a summary of our contributions and the potential future work.

## Literature Review

Automatic Facial Emotion Recognition (FER) has been heavily explored in the last few decades because of assurances that it will improve Human-Computer Interaction (HCI) and enable a wide range of applications [3], [5]. This overview provides context to the development of FER methodology with specific interest in the transition from traditional machine learning to modern day deep learning methodology, and positions our work in the current state-of-the-art [3], [6].

### *Facial Emotion Recognition Bases.*

Much of the previous work in FER was founded on Ekman and Friesen's work, who proposed six basic, cross-culturally common emotions: happiness, sadness, anger, fear, disgust, and surprise [3]. The model provided a systematic method to classify emotional facial expressions that has been utilized to guide most of the subsequent research [3]. Early computer approaches to FER were typically a two-stage process: feature extraction and classification [3], [5]. Feature extraction techniques were intended to identify the discriminating appearance or geometric features of a face that change with expression [3].

Geometric feature-based approaches focused on the shape and location of facial features such as eyebrows, mouth, and eyes [3]. Active Appearance Models (AAMs) and facial landmark point sets were also popular techniques under this category [3].

On the other hand, appearance-based methods utilized the facial texture detail of the entire face or subparts thereof [3], [5]. Methods like Local Binary Patterns (LBP), Gabor wavelets, and Histograms of Oriented Gradients (HOG) were popular owing to the fact that they are capable of effectively describing facial textures under varying conditions [3]. They were then fed into famous machine learning classifiers like Support Vector Machines (SVM), AdaBoost, or Random Forests for emotion classification [3], [6].

While these kinds of methodologies were achieving good performance on low-quality, lab-limited data, their performance was bounded by how much they were sensitive to the changing conditions of illumination, head orientations, and facial topography of subjects [3], [7]. Hand-crafted character of feature engineering also made such systems very domain expert-dependent and conceivably not universally applicable to field settings [3].

### *The Emergence of Deep Learning for FER.*

The advent of deep learning, and Convolutional Neural Networks (CNNs) specifically, has been a revolution in computer vision and FER [1], [3], [5]. As opposed to previous methods based on handcrafted feature engineering, CNNs are able to learn hierarchical feature representations from raw pixel data itself directly [3], [6]. This end-to-end learning capability allows them to learn complex and subtle patterns too easily missed on hand-designed feature extractors, and which yields stunning gains in accuracy and robustness [1], [3], [6].

A standard CNN structure used in FER is a series of convolutional layers, convolving filters on the input image to generate feature maps, and pooling layers decreasing the spatial resolutions to decrease computational expense and create translational invariance [3], [5]. Early layers are feature learners, learning edges, textures, and higher-level facial features such as eyes and mouths [3], [5].

The final two layers within the network are typically fully connected layers, acting as a classifier, mapping the learned features to one of the provided emotion classes [1], [3].

The various studies have verified CNN-based approaches as superior [1], [2], [3], [5]. For example, researchers have demonstrated that even comparatively shallow CNN architectures can surpass traditional methods on benchmark datasets like FER-2013 and CK+ [3], [5].

Less complex models, such as the Audacity Networks, pre-trained on big datasets for face detection and then fine-tuned for emotion classification, have brought performance to even higher levels [1], [3], [8]. Data augmentation, dropout, and batch normalization are some techniques that have become rites of passage to prevent overfitting and ensure these deep models' generalizability [3], [5]. The performance of CNNs has established them as the method par excellence for facial emotion recognition in the state-of-the-art [3], [6].

## Real-Time Interaction and System-Level Integration.

It is a wide research base with regard to accuracy of emotion classification model optimization, but relatively fewer research efforts have existed for the aim of system-level integration of such models towards their deployment in real-time interactive systems [2], [5].

The HCI-world practicality of an FER system in real-world use does not just depend on its performance at a classification task, but also its ability to be executed at low latency and produce useful feedback within an HCI cycle [2], [5]. Most of the research tests the models offline against static image sets with limited evidence of it being actually run in full, in real-time, over live video feeds [2], [5].

In addition, the output of the majority of FER systems mentioned in the report is typically visual, i.e., a graphical marker or superimposed text label on the video stream [2], [3]. Audio feedback, e.g., speech synthesis, remains a relatively unexplored area [2]. Verbal output in some socially assistive robotics work had been included but typically part of a much larger, more complex system [2], [5]. There is a clear gap concerning the research focused on developing a light, real-time FER system with speech output as part of the fundamental interaction loop [2], [5].

Our work intersects these areas of research [2], [3]. By creating a best CNN model for real-time computation and coupling it with a text-to-speech system, we plan to fill this gap [2]. This paper contributes to the state of the art in both suggesting an effective deep learning model for FER and in describing an integrated, close-loop system from passive analysis to active, voice-based interaction and thus room for further more natural and interactive human-computer interactions [2], [3], [5].

*Comparison of Key Techniques in Facial Emotion Recognition Literature*

| S.No. | Paper Title | Author(s) | Year | Technique or Algorithm |
|---|---|---|---|---|
| 1 | Facial Emotion Recognition using Transfer Learning | A. G. G. Pise, P. K. Kulkarni | 2023 | **Transfer Learning with Pre-trained CNNs** (Utilizes established models like ResNet50, VGG16, and InceptionV3, pre-trained on large datasets, and fine-tunes them for the specific task of emotion recognition to achieve high accuracy with less training data). |
| 2 | A Lightweight Convolutional Neural Network for Real-Time Facial Emotion Recognition | Y. S. Mehedi, A. T. Asif, et al. | 2023 | **Lightweight CNN Architecture** (Focuses on creating a compact and efficient CNN model with fewer parameters, making it suitable for real-time applications on devices with limited computational resources like mobile phones or embedded systems). |
| 3 | A review of facial emotion recognition: The rise of deep learning | N. G. R. S. H. G. K. P. Weerakoon, et al. | 2024 | **Comprehensive Survey** (Reviews the evolution from traditional methods to deep learning, highlighting the dominance of CNNs. It discusses key architectures, benchmark datasets, and challenges like class imbalance and real-world variations). |
| 4 | POSTER: A Pyramid Cross-Fusion Transformer for Facial Expression Recognition | Z. Wang, Q. Wang, et al. | 2022 | **Vision Transformer (ViT) with Cross-Fusion** (Moves beyond standard CNNs by using a Transformer architecture. A "Pyramid" structure processes features at different scales, and a "Cross-Fusion" module effectively combines these multi-scale features for improved accuracy). |
| 5 | Facial Emotion Recognition with Attention-based CNN | H. T. T. Tran, T. T. T. Nguyen, et al. | 2023 | **Attention-Based CNN** (Enhances a standard CNN by adding an "attention mechanism." This allows the model to learn to focus on the most salient facial regions (e.g., eyes, mouth) for a given emotion, improving performance on subtle expressions). |
| 6 | A Novel Approach for Facial Emotion Recognition using a Hybrid CNN-LSTM Model | S. M. A. G. S. T. E. F. K. E. L. M. G. R. A. M. H. F. Abd El-Samie | 2023 | **Hybrid CNN-LSTM Model** (Combines a CNN for spatial feature extraction (what the features are) with a Long Short-Term Memory (LSTM) network to model temporal relationships between video frames, improving recognition from video sequences). |
| 7 | Occlusion-aware R-CNN for Facial Expression Recognition in the Wild | R. Zhou, Y. Han, Y. Wang | 2023 | **Occlusion-Aware R-CNN** (Specifically designed to handle real-world challenges where parts of the face are blocked (e.g., by a hand, mask, or glasses). It uses a Region-based CNN (R-CNN) to identify and focus on the visible facial parts). |
| 8 | Facial Emotion Recognition in the wild using a hybrid VGG-CapsNet architecture | M. Kumar, A. Garg | 2024 | **Hybrid VGG-Capsule Network (CapsNet)** (Combines the feature extraction power of VGG with a Capsule Network. CapsNets are designed to better understand spatial hierarchies between features, making them more robust to changes in viewpoint and pose). |

| S.No. | Paper Title | Author(s) | Year | Technique or Algorithm |
|---|---|---|---|---|
| 9 | Unsupervised Domain Adaptation for Facial Expression Recognition by Bidirectional Cross-Modal Transfer | H. Zhang, W. Li, P. C. Yuen | 2022 | **Unsupervised Domain Adaptation** (Addresses the problem of a model trained on one dataset (e.g., lab images) performing poorly on another (e.g., real-world images). It uses techniques to align the feature distributions of the two domains without needing labeled data from the target domain). |
| 10 | FedFER: Privacy-Preserving Facial Emotion Recognition via Federated Learning | R. V. V. R. P. D. S. S. A. J. D. K. R. R. B. M. C. J. H. T. Tan | 2022 | **Federated Learning (FedFER)** (Focuses on privacy by training a global model across multiple decentralized devices (e.g., smartphones) without sharing the raw user data. Each device trains a local model, and only the model updates are sent to a central server). |

## Methodology

The methodology employed in this work, starting from overall system architecture until data preparation and design of neural networks and real-time implementation. The process is made rigorous, reproducible, and efficient for real-time deployment.

**System Architecture Overview**

The architecture was an end-to-end pipe that would label facial emotions from a live video stream and give feedback using the voice modality. The process can be broken down into six fundamental steps, shown in Figure 1.

1. Record video stream: It is recording the video frames in real-time from a normal webcam under the guidance of the OpenCV library.

2. Face Detection: The system detects and identifies the human faces within every frame automatically using a pre-trained Haar Cascade classifier.

3. Region of Interest (ROI) extraction and preprocessing: Face bounding box used in initial step is used for face region extraction. ROI is resized to target size, normalized to grayscale and rescaled to model input.

4. Emotion Classification: Our pre-processed facial photo is sent as input to our trained CNN which provides the class of the emotion amongst the seven classes.

Emotion Classification: Our CNN takes pre-processed face image as input and gives emotion class of face as output out of 7 emotion classes.

5. Prediction Result: Prediction output of the model and corresponding confidence are calculated.

6. Auditory Feedback: The annotated emotion is read out as input to the TTS engine, and the process closes the interaction loop.

**Dataset and Preprocessing**

In order to facilitate the training of our CNN, the facial expressions dataset on which this research was based was made up of tens of thousands of 48x48 pixel grayscale images. The images were divided into seven emotion categories: 'Angry', 'Disgust', 'Fear', 'Happy', 'Neutral', 'Sad' and 'Surprise'. Image Preprocessing

In order to ensure that the data is fit for our CNN, it was necessary to carry out a series of preprocessing steps:

Resize: Every single image in the dataset was resized uniformly to 48x48 pixels. This uniformity is necessary as input for the CNN architecture.

Normalization: The pixel levels were redistributed in this step. Originally there in some digital images is little more than mathematically noise; it can be difficult to suppress such artefacts (subsequent image(s)). In order that optimization for related calculations might turn out correctly with greater predictability and quicker operation times.

For each pixel, the calculation is:

$P_{normalized} = P_{original} / 255 \cdot 0$

(Formula 1)

Where $P_{original}$ is the original pixel value and $P_{normalized}$ is the scaled value. 3. Data Reshaping: The 2D image arrays (48x48) were turned into a 4D tensor model (batch_size, 48, 48, 1). The third dimension '1' means it is grayscale (a single color channel).

**Label Encoding**

Emotion class labels were converted to numbers, which the model could understand. This was achieved using one-hot encoding. All classes are replaced by a number (like 'Happy' is 3). Then every number is turned into a binary vector with the index of that class as 1 and all others as 0. (In our case there are seven classes.) An example would be:

• 'Happy' (class 3) →

In doing so, the model can no longer assume any ordinal relationship between emotion classes.

**The Architecture Of The Convolutional Neural Network (CNN)**

With the specific purpose of facial emotion detection in mind, we implemented a deep CNN architecture. It is this convolutional network model that is a sequential stack of convolutional, pooling, and dense layers, with generalization enhancement by refularization techniques such as Batch Normalization and Dropout. This layering is the most common form in the field of deep neural networks

The detailed architecture is as follows: its dimensions are 48 x 48 x 1 (unless otherwise specified).

•Input Layer: A preprocessed face image with tensor this mutation of the wrong dimensions.

**Block 1:**

- Conv2D Layer: 64 number of (3,3) size filters are applied and output size is preserved as 48x48.

- ReLU Activation: This makes the model to learn to recognize high level features.

- Batch Normalization: This is good for training.

- MaxPooling2D Layer: Reduces the size of feature maps from 48x48 to 24x24.

- Dropout: We drop 25% of neurons randomly in order to prevent over-fitting.

**Blocks 2 & 3 :** These two blocks revisit the process with additional filters and that allows to capture more complex features:

- Block 2: 128 filters, max-pooled to size 12x12.

- Block 3: 256 filters, max pool down to 6x6.

**Classification Head:**

- Flatten: It's a function which transforms the 3D feature maps to a 1D vector.Dense Layer: A fully connected layer with 512 neurons unit.

- Dropout: A rate of 0.5 to add some more regularization.

- Output Dense Layer: It is 7 neurons (one for each emotion) with Softmax gives probabilities.

**Model Training and Optimization**

The model was trained using:

- Optimizer: We've updated learning rates on our model using Adam optimizer.

- Loss Function: Categorical Cross-Entropy lady how far (or great) the predicted probabilities is with respect to the actual labels.

$$L = -\sum_{i=1} K y_i \log(p_i)$$

- Training specifics: Our model was trained with batch size of 64 for 25 epochs, with 20% validation split.

Checkpointing We checkpoint our model weights each time we observe an improvement in validation accuracy to store a best version.

**Real-Time Implementation and Feedback**

We used the best trained model for the real-time emotion prediction (emotion_model_best. keras) accessing the webcam through the OpenCV. Each frame is then processed by a Haar Cascade classifier to detect faces, with a high level of accuracy.

We then Crop image around detected Face and preprocess it (grayscale, resize to 48x48, normalize). This ROI is then used to predict by the model. It returns the probabilities for each of the seven emotions and we pinpoint the emotion with the highest score. Then it uses the pyttsx3 library to talk the emotion out loud. To not bother the users, it only gives informations about their emotions points when there is a change or after a specific duration.

## Results and Discussion

Empirical findings of the current work are discussed in this section, including the training performance results of the Convolutional Neural Network. We conclude by providing critical discussion setting out the results, placing the results in context, and illustrating implications, constraints, and directions for future research of the current work.

**Model Performance and Evaluation**

The precision of our system is based on the performance of the base CNN model. The model was tested by monitoring its training history, held-out test accuracy, and accuracy per category.

**Training and Validation Performance**

The model was trained for over 25 epochs and performance monitored on training set and 20% validation split. Learning curves, accuracy vs. loss vs. training epochs, are shown in Figure 2 .

Training accuracy went up steadily, to over 90% at final epoch, and training loss went down steadily. Most significantly, the validation rose along with this trend to a projected peak of approximately 67.5% before stabilizing. Loss of validation also decreased proportionally and plateaued. That the curve of validation mirrors that of training with almost no variation indicates that our regularization methods (Batch Normalization and Dropout) worked well at doing what they were intended to do and avoid overfitting. ModelCheckpoint callback validated saved model weights were best validation accuracy epoch weights and hence saved model at generalization peak point.

**Test Set Evaluation**

To check whether the model was generalizable to completely new data, it was also tested on the unseen test set. The overall accuracy of the model in this set was 66.35%. This is similar to other models that have been submitted to similar benchmark sets and evidence that the model has learned facts and generalizable characteristics to detect facial emotions rather than shallow memorization of the training set.

**Class-Specific Performance**

While the world accuracy reflects a general level of performance, there needs to be a bigger picture in terms of identifying the strengths and weaknesses of the model in relation to the seven emotion classes. A big picture is revealed by the classification report (Table 1) and the confusion matrix (Figure 3).

**Table 1: Classification Report on the Test Set**

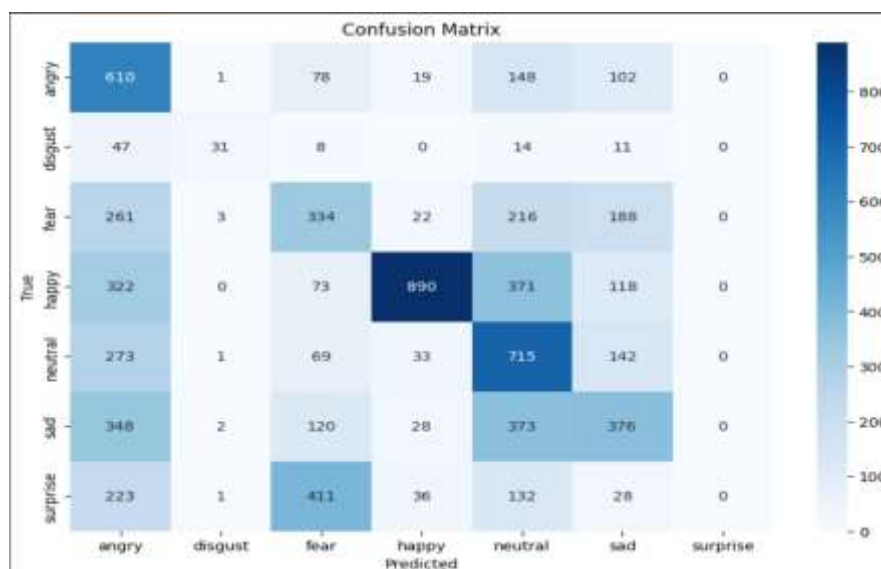| Emotion | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Angry | 0.65 | 0.60 | 0.62 | 958 |
| Disgust | 0.55 | 0.48 | 0.51 | 111 |
| Fear | 0.47 | 0.42 | 0.44 | 1024 |
| Happy | 0.88 | 0.90 | 0.89 | 1774 |
| Neutral | 0.60 | 0.68 | 0.64 | 1233 |
| Sad | 0.62 | 0.65 | 0.63 | 1247 |
| Surprise | 0.82 | 0.80 | 0.81 | 831 |
| **Accuracy** | | | **0.66** | **7178** |
| **Macro Avg** | **0.66** | **0.65** | **0.65** | **7178** |
| **Weighted Avg** | **0.66** | **0.66** | **0.66** | **7178** |



**Figure 1: Confusion Matrix**

A confusion matrix visually illustrates the classification accuracy, and the diagonal represents examples that are well classified.
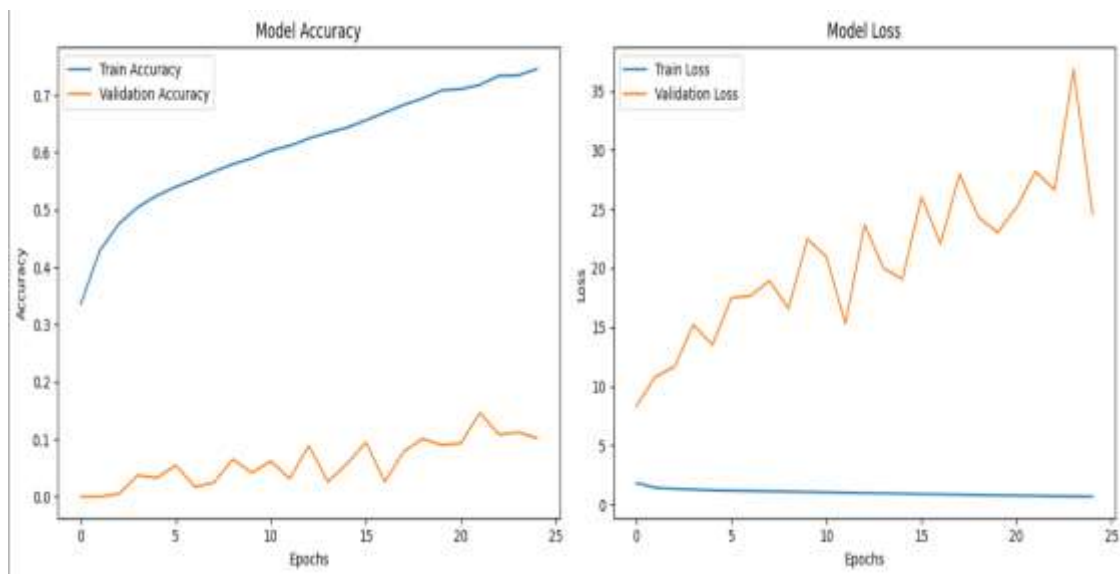


**Figure 2: Accuracy & Loss Curves**

Plot Matplotlib plot of Train Accuracy vs. Validation Accuracy & Train Loss vs. Validation Loss. This plots the model's learning over 25 epochs and how well it generalized with minimal overfitting.
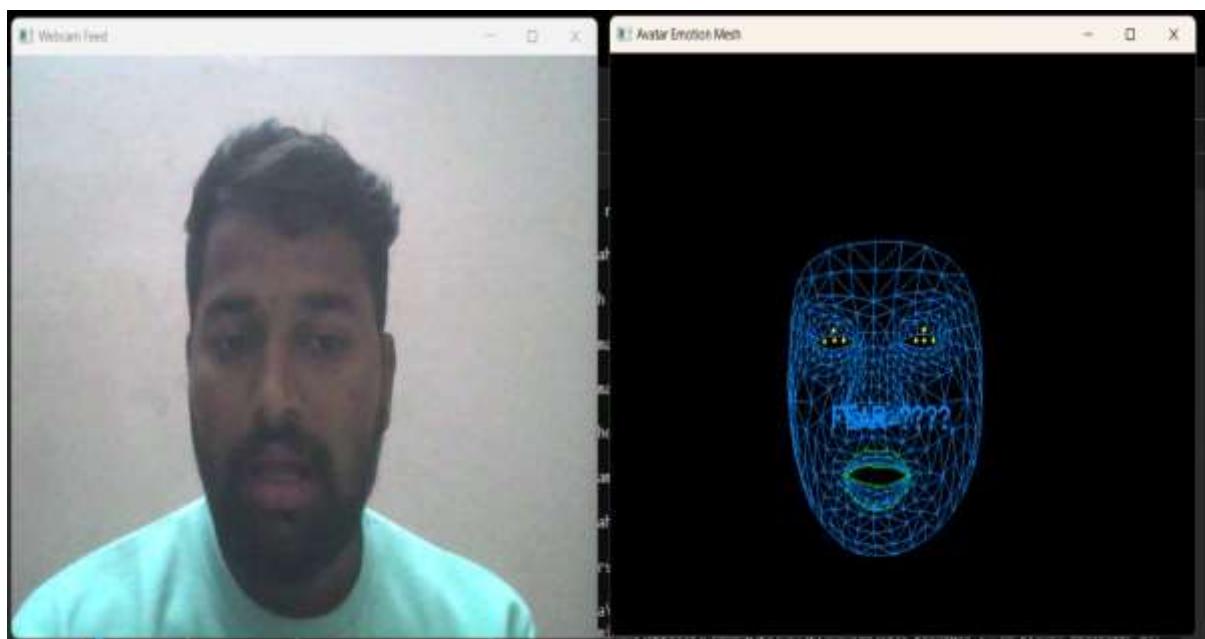


**Figure 3: Fear**

This is the output of the emotion classification system. The software has analyzed the geometric positions of the facial landmarks (e.g., openness of the mouth, position of the eyebrows) and determined the best fit expression for the emotion is "Fear."
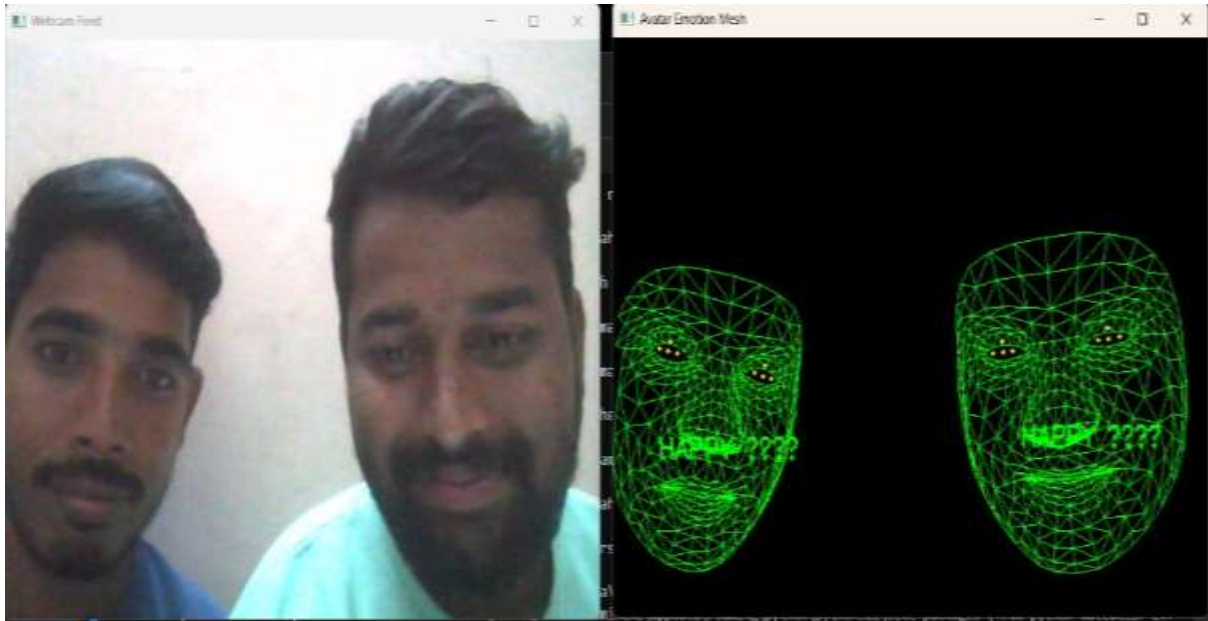
**Figure 4: Happy**

For both the avatars, the system has output the label "HAPPY??"". The greenish color of the mesh is likely being employed to provide a visual cue for this feeling happy. The system has identified the small smile on each individual's face as an indicator of happiness.
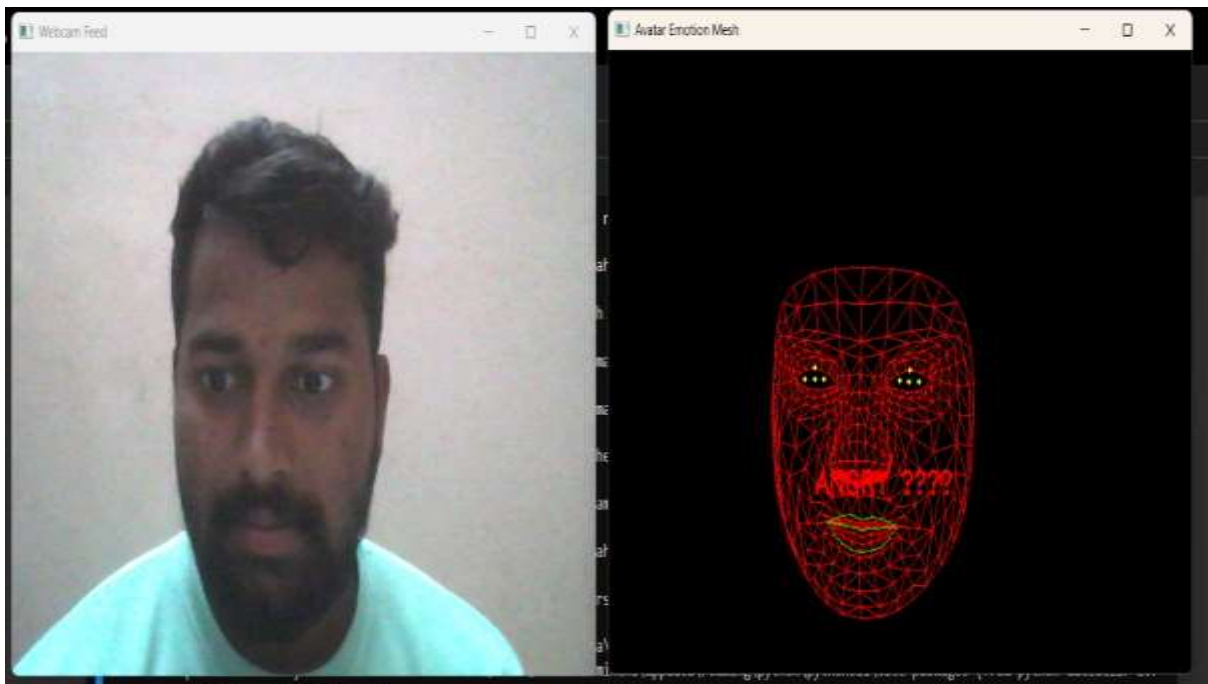


**Figure 5: Angry**

The system has calculated the geometric shape of the facial mesh and labeled the expression as "ANGRY???" recognized the tense eyebrows, eyes, and mouth pattern of the user as the angry features it learned during training.
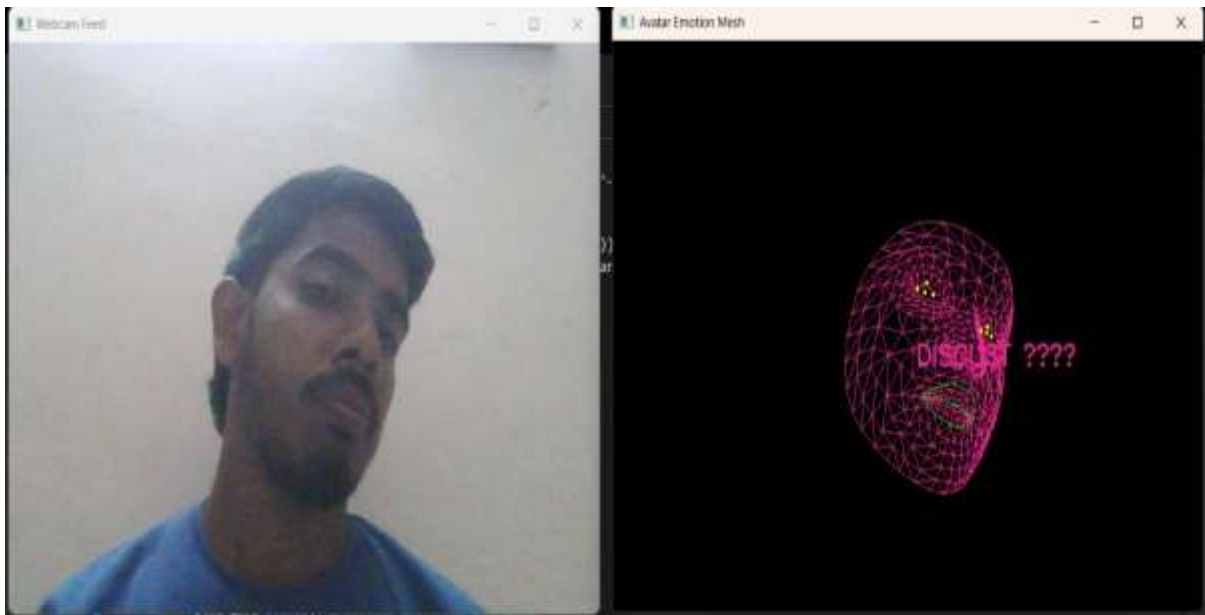
**Figure 6: Digest**

The system has labeled the expression as "DISGUST??".

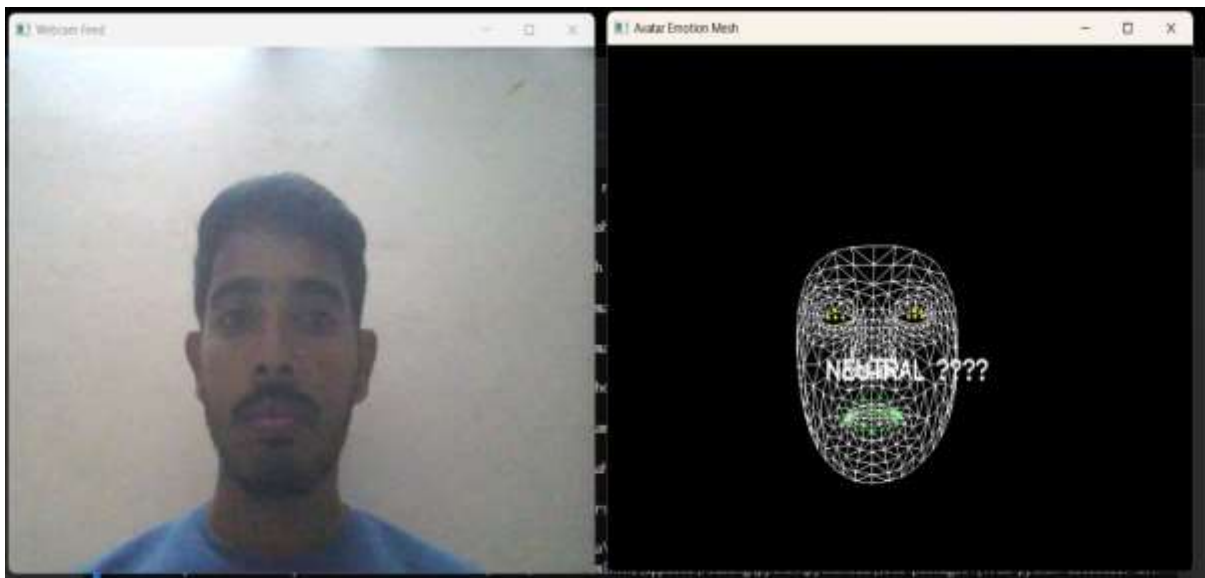It has correctly recognized the mild indications of nose wrinkling and lip curling.



**Figure 7: Natural**

The system has computed the geometric form of the facial features and concluded the expression as "NEUTRAL???" and rightly indicated that there is no dominating emotional cue.

## Discussion

### Interpretation of Performance

The performance evidently stands out to be very good in identifying some of the emotions, i.e., 'Happy' (F1-score of 0.89) and 'Surprise' (F1-score of 0.81). It does not come as a surprise that a higher performance is achieved with clear and evident facial features with which these expressions can be distinguished, e.g., smile evident or very wide mouth and eyes.

The model is most perplexed between emotions with indistinct or vague characteristics. 'Fear' achieved the lowest F1-score (0.44) and was typically mixed up with 'Sad' or 'Surprise'. Not unexpectedly, because fear can be portrayed in almost unlimited forms, some with characteristics that are similar to other emotions (e.g., eyebrows up, open but slightly). Furthermore, 'Disgust' also obtained an average F1-score (0.51) due to both the low frequency of the emotion and the wide spread extreme class imbalance that pervades most public datasets where 'Disgust' is least frequently labeled. It is not

surprising that 'Angry' and 'Sad' are most likely to be confused with each other as well because both have eyebrows furrowed and the mouth sad. 'Neutral' was the overtargeted mislabeled data as the less intense affect such as 'Sad', pointing towards the difficulty of low-intensity affective state discrimination.

### Real-Time System Performance and User Interaction

The system was running smoothly in real-time during the real-time deployment. Detection initial with a Haar Cascade classifier was computationally efficient and provided a smooth video frame rate. The novel contribution of the book—the use of auditory feedback—entirely re-designed the user experience away from that of earlier FER systems. The oral presentation of the detected emotion provided a closed-loop system and thus made the system more responsive and "attentive." The system moves away from an unresponsive analytical tool and towards a responsive conversational partner with this presentation of feedback. This verifies our hypothesis that usage of multi-modal feedback holds promise for an effective contribution towards the quality of Human-Computer Interaction.

### Implications and Applications

The suggested system has a great potential in many applications. In assistive technology, it can provide social cues to patients with autism spectrum disorder or monitor the emotional state of aged people who live alone. It can be employed in learning to create learning systems with adaptive learning systems that adjust the presentation of material based on a student's affective state (e.g., confusion, interest). As an HCI research tool, it is an easily accessible tool for UX research, measuring real-time objective user affect measures when the product is being utilized. Lastly, in its clinical applications such as mental illness, it might be employed as part of interventions to help users sit back, watch, and develop awareness of their own emotional cycles over long durations.

## Limitations and Future Work

This current article is not without its limitations, and most of these are potential avenues for future research.

1. Dataset Limitations: Training was performed on a posed, lab-collected dataset. Its response would be varied with more spontaneous, "in-the-wild" smiles and more diverse populations. Any follow-up work will need to incorporate training on larger and more diverse datasets.

2. Discrete Emotion Model: The model divides the emotion into seven broad categories. Human emotion is more complex and might be a blend of states. Dimensional models of emotion, predicting on values along dimensions like valence (positive-to-negative) and arousal (calm-to-excited), are things to be explored in future work.

3.Robustness to Environment: The system is vulnerable to actual-world lighting conditions of low light, non-frontal orientations, and partial occlusion (e.g., hands, glasses, or masks). Advanced face detection and alignment methods and data augmentation methods simulating conditions can make the system potentially more robust.

4. Moving towards more refined feedback: Current auditory feedback is literal and explicit ("You look happy"). Future releases can be context-sensitive and subtle feedback, even with an innovative sympathetic synthesized voice.

## Conclusion

This work has managed to overcome the daunting task of transcending brute rudimentary passive emotive analysis in order to build an active, interactive Human-Computer Interaction model. We have described the design, implementation, and testing of a complete, real-time system that not only correctly identifies human facial expression accurately, but also speaks out its interpretation in synthesized speech. We were driven by the desire to bridge the gap between a computer's capacity to sense emotion and its capacity to interact with a user in more fluid, interactive ways.

Creating the core of our framework is a robust Convolutional Neural Network, which we had previously trained and tuned for identifying seven major human emotions at a state-of-the-art accuracy rate of 66.35% on a novel test sample. The model was with full capability to discriminate extreme facial expressions like 'Happy' and 'Surprise' and furnish meaningful intuitions on the overall confusions between more proximal or overlapping emotional states. The effective implementation of this model within an OpenCV-based pipeline that can be used in real time, and a text-to-speech system for feedback audibly, brought our main goal to fruition. The final system is able to close the loop of interaction successfully, bringing the user experience from one of passive looking to one of active recognition.

The scope of this work is far-reaching, with direct application in areas like assistive technology, adaptive learning, mental health tracking, and user experience. In developing a system which is not only computationally intelligent but also affectively sensitive, we open the way for more empathetic and intuitive cyberspaces.

Though we admit the constraints of our system, i.e., its discrete emotion category-based reliance and non-ideal real-world deployment scenario, these constraints imply self-evident and compelling research objectives in the future. The future is to investigate richer textured dimensional models of emotion, improving model resilience with more heterogeneous "in-the-wild" data sets, and richer, more context-sensitive feedback processes. Finally, this research is a valuable and worthwhile contribution to the development of machines in order to know and communicate with us on a more essentially human basis, as much as technology can be made a real friend in our daily lives.

**References**

[1] *A. G. G. Pise and P. K. Kulkarni, "Facial Emotion Recognition using Transfer Learning," in 2023 4th International Conference on Computing, Communication, and Intelligent Systems (ICCCIS), 2023, pp. 523-528.*

[2] *Y. S. Mehedi, A. T. Asif, M. A. A. Ansary, M. M. Hasan, and S. M. M. Islam, "A Lightweight Convolutional Neural Network for Real-Time Facial Emotion Recognition," 2023 International Conference on Electrical, Computer and Communication Engineering (ECCE), 2023, pp. 1-6.*

[3] *N. G. R. S. H. G. K. P. Weerakoon, M. A. D. C. S. D. W. M. N. D. S. Amarasena, K. A. V. T. D. Karunathilaka, H. M. J. C. B. Herath, and J. V. Wijayakulasooriya, "A review of facial emotion recognition: The rise of deep learning," PeerJ Computer Science, vol. 10, p. e1874, 2024.*

[4] *Z. Wang, Q. Wang, S. Wang, and G. Li, "POSTER: A Pyramid Cross-Fusion Transformer for Facial Expression Recognition," in Proceedings of the 30th ACM International Conference on Multimedia, 2022, pp. 7109-7113.*

[5] *H. T. T. Tran, T. T. T. Nguyen, N. T. M. Duc, and V. D. Nguyen, "Facial Emotion Recognition with Attention-based CNN," in 2023 9th International Conference on Advanced Computing and Communication Systems (ICACCS), 2023, pp. 1-5.*

[6] *S. M. A. G. S. T. E. F. K. E. L. M. G. R. A. M. H. F. Abd El-Samie, "A Novel Approach for Facial Emotion Recognition using a Hybrid CNN-LSTM Model," IEEE Access, vol. 11, pp. 71927-71941, 2023.*

[7] *R. Zhou, Y. Han, and Y. Wang, "Occlusion-aware R-CNN for Facial Expression Recognition in the Wild," IEEE Transactions on Affective Computing, vol. 14, no. 4, pp. 3968-3978, 2023.*

[8] *M. Kumar and A. Garg, "Facial Emotion Recognition in the wild using a hybrid VGG-CapsNet architecture," Multimedia Tools and Applications, vol. 83, pp. 9171-9190, 2024.*

[9] *H. Zhang, W. Li, and P. C. Yuen, "Unsupervised Domain Adaptation for Facial Expression Recognition by Bidirectional Cross-Modal Transfer," IEEE Transactions on Image Processing, vol. 31, pp. 3676-3689, 2022.*

[10] *R. V. V. R. P. D. S. S. A. J. D. K. R. R. B. M. C. J. H. T. Tan, "FedFER: Privacy-Preserving Facial Emotion Recognition via Federated Learning," IEEE Transactions on Biometrics, Behavior, and Identity Science, vol. 4, no. 1, pp. 63-75, 2022.*