# Emotion Detection in Facial Expressions and Speech Using Deep Hybrid Learning

## *Mr. Pragash K[1], Deepak R[2], Gopinath R[3], Shasteeswaran Ravianand Tharmiya[4], Rohit S. K[5]*

[1]Associate Professor, Department of Artificial Intelligence and Data Science, Sri Manakula Vinayagar Engineering College, Puducherry, India
pragashkaliyan@gmail.com
[2,3,4,5]Bachelor of Technology – Artificial Intelligence & Data Science, Sri Manakula Vinayagar Engineering College, Puducherry, India,
[2]dr194304@gmail.com, [3]vijayagopinath0112@gmail.com, [4]trshasteeswaran2003@gmail.com, [5]rohitsk2003@gmail.com

## ABSTRACT

This project presents the development of an automatic emotion detection system that integrates facial expression analysis and speech processing through deep hybrid learning. By advancing emotion recognition with both visual and audio inputs, the system aims for high precision in detecting human emotions. Facial expression data is gathered from the FER-2013 dataset, while emotional speech data is sourced from the RAVDESS dataset. Convolutional neural networks (CNNs) extract features from facial images, and recurrent neural networks (RNNs) with LSTM units capture temporal dependencies in speech. The fusion of these modalities provides a comprehensive representation of emotional states, ensuring both visual and auditory information contribute to detection accuracy. Traditional machine learning algorithms like Support Vector Machines (SVM) and Random Forests are employed for classification, with categorical cross-entropy as the loss function. To evaluate performance, accuracy, precision, recall, and a confusion matrix are utilized, validating robustness across multiple emotion categories. This multimodal approach outperforms

Single-modality methods, achieving. Significantly improved accuracy and reliability in emotion detection.

Overall, the project demonstrates the effectiveness of combining facial expression analysis with speech processing, offering a more human-like, context-aware emotion detection system.

## 1. INTRODUCTION

Understanding emotions is central to the design of many modern intelligent systems. Emotions drive human interactions, influencing how we communicate, make decisions, and connect with others. From virtual assistants that aim to make conversations more engaging, to mental health applications that monitor emotional well-being, and customer service platforms that tailor responses to customer moods, emotion detection technology has become increasingly essential. The goal of emotion detection is not simply to identify a specific emotion but to enable systems to respond more naturally and empathetically, making interactions with technology feel more human and context-aware.

Emotion detection is a complex task because emotions are multifaceted and can be expressed in various ways. Humans communicate emotions through facial expressions, voice tone, body language, and even subtle cues like pauses or changes in breathing. For example, a person expressing frustration might show a furrowed brow and clenched jaw, speak in a faster, louder tone, or exhibit tense body language. Each of these cues contributes valuable information to accurately interpret the person's emotional state. However, capturing and processing these nuances is a challenge, especially for traditional emotion detection systems that typically rely on just one type of input.

Single-modality models, which focus on either facial expressions or speech alone, have been the backbone of earlier emotion detection efforts. Facial expression recognition systems, for instance, use image data to analyze changes in facial features, such as the movement of eyebrows or the shape of the mouth, to infer emotions. Similarly, speech-based emotion detection relies on audio data to capture elements like pitch, volume, and rhythm. Although these single-modality models can be effective in controlled settings, they often struggle in real-life applications. For example, facial expression recognition may fail in low-light conditions or with individuals who naturally express emotions differently due to cultural differences. Likewise, speech-based models may misinterpret emotions in noisy environments or with speakers who have unique vocal patterns.

The limitations of these older, single-modality approaches reveal the need for more advanced models capable of integrating multiple types of data. By combining different forms of input—such as both visual and auditory cues—emotion detection systems can achieve a more holistic and accurate understanding of human emotions. This hybrid, or multimodal, approach allows systems to compensate for limitations in one modality by drawing on

information from another. In essence, when facial expressions alone are ambiguous, vocal tone and rhythm can provide critical context, and vice versa, creating a fuller and more reliable emotional profile.

### 1.1. Need for a Hybrid Model

A hybrid learning approach that combines CNN and LSTM models is especially promising for advancing emotion detection. By looking at both facial expressions and vocal cues simultaneously, these hybrid models provide a richer, more reliable understanding of emotions, even in more complex scenarios. For instance, when facial expressions might be difficult to read, voice tone can fill in critical context and vice versa. This dual-modality approach doesn't just boost classification accuracy; it also makes the system more capable of interpreting subtle, context-dependent emotions.

Hybrid models tackle some of the common shortcomings of single-modality systems through multimodal fusion, where the outputs from CNN-based facial analysis and LSTM-based speech analysis combine into a unified emotional profile. This way, the system benefits from the strengths of both facial and vocal inputs, making it better suited for diverse and dynamic real-world applications, from customer service to health monitoring and interactive learning platforms.

## 2. LITERATURE SURVEY

Research into emotion detection has rapidly advanced, thanks to the rise of deep learning and multimodal processing techniques. This literature review examines key developments in emotion detection across several relevant areas, including deep learning for facial expression and speech analysis, hybrid approaches, and the integration of multimodal systems.

### 2.1. Emotion Detection using Deep Learning

Deep learning has significantly enhanced the ability to interpret complex data patterns, particularly in fields like computer vision and speech processing. Convolutional Neural Networks (CNNs) are widely used in facial expression analysis due to their capacity to extract intricate spatial features. CNNs can identify subtle changes in facial regions—such as the eyes, eyebrows, and mouth—that reveal emotional cues. Notable advancements in CNN-based emotion detection include models trained on established datasets like FER-2013, which have achieved promising accuracy levels. These models can generalize well in controlled environments but face limitations in diverse, real-world scenarios where lighting, angles, and cultural differences in expression can affect accuracy. Parallel to CNN advancements, Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, have proven effective for speech-based emotion detection. LSTMs excel at capturing sequential data patterns, such as changes in pitch, tone, and rhythm in speech, which are essential indicators of emotional state. Studies leveraging LSTMs, often with datasets like the RAVDESS (Ryerson Audio-Visual Database of Emotional Speech and Song), show that these models can differentiate between emotions such as happiness, sadness, anger, and surprise based on vocal features. However, similar to facial models, these speech models also face challenges in real-world contexts, including handling background noise and speaker variability.

### 2.2. Hybrid Emotion Detection Models

In response to the limitations of single-modality models, researchers have explored hybrid approaches that combine CNNs and LSTMs for a more comprehensive emotion detection system. Hybrid models leverage the strengths of each modality, allowing for both spatial analysis of facial features and temporal analysis of vocal cues. Studies have shown that hybrid models can outperform single-modality systems by providing a fuller emotional context, especially in cases where one cue may be ambiguous. For instance, a study by Zeng et al. (2021) highlighted that hybrid models achieved higher classification accuracy in detecting complex emotions compared to CNN or LSTM models alone, underscoring the effectiveness of a multimodal approach in emotion recognition.

### 2.3. Multimodal Fusion Techniques

Fusion techniques play a crucial role in hybrid emotion detection models, as they determine how data from multiple modalities is combined to produce a final prediction. Multimodal fusion approaches are typically categorized into three types: early fusion, late fusion, and intermediate fusion. Early fusion involves combining raw data from both modalities before feature extraction, which can improve integration but requires significant computational power. Late fusion, on the other hand, combines the outputs of each modality's separate processing pipeline, which can reduce computational load but may miss out on subtle intermodal relationships. Intermediate fusion, which combines data at the feature extraction stage, is often considered a balanced approach, allowing the model to capture relevant features from both modalities while maintaining efficient processing. Recent studies have emphasized the potential of attention mechanisms in multimodal fusion, as they allow the model to focus on the most relevant features in each modality. For example, Zhao et al. (2020) integrated attention layers into a hybrid model, which selectively focused on critical vocal or facial features based on the context, significantly enhancing classification accuracy. These findings highlight how sophisticated fusion techniques, especially when incorporating attention mechanisms, can improve the interpretability and accuracy of emotion detection models.

*2.4. Applications of Emotion Detection*

The applications of emotion detection span various industries, including healthcare, customer service, education, and entertainment. In mental health, for example, emotion detection can assist therapists by monitoring patients' emotional responses, allowing for more personalized treatment. In customer service, emotion detection models integrated with chatbots and virtual assistants help adjust responses to match customer sentiment, improving user satisfaction. Educational tools that utilize emotion detection can gauge student engagement, offering real-time adjustments to keep learners motivated. Entertainment platforms, especially virtual reality and gaming, also benefit from emotion detection by creating interactive, adaptive environments that respond to players' emotions, enhancing the immersive experience. Research consistently shows that these applications have higher success rates and user engagement when emotion detection systems incorporate multimodal, hybrid approaches, highlighting the need for further development in this area.

*2.5. Challenges and Future Directions*

Despite recent advances, several challenges remain in the development and deployment of emotion detection systems. Multimodal models require large amounts of diverse training data to generalize well across different demographic groups and environmental settings. Ensuring privacy and data security, particularly in sensitive applications like mental health, is another significant concern. Additionally, the computational demands of multimodal systems can be a barrier to real-time processing, especially in mobile and edge computing applications. Future research may focus on improving data efficiency through techniques like transfer learning, which could enable emotion detection models to learn from smaller datasets. Additionally, advances in lightweight model architectures, such as mobile-optimized CNNs and LSTMs, could make emotion detection more accessible in real-time applications on lower-power devices. Addressing these challenges will be key to creating more reliable, scalable, and ethical emotion detection systems for practical, real-world use.

| Aspect | Old Ideas | New Ideas |
|---|---|---|
| Facial Emotion Recognition | Manual coding (FACS), CNNs for static features | CNN + LSTM hybrid for dynamic, real-time emotion tracking |
| Speech Emotion Recognition | DNN for speech data only | CNN + LSTM with advanced speech preprocessing (noise reduction, pitch normalization) |
| Multimodal Emotion Detection | Early/Late fusion of modalities | Intermediate fusion with attention mechanisms to weigh data importance dynamically |
| Data Synchronization | Basic fusion without synchronization | Intermediate fusion with attention to address timing issues between speech and facial data |
| Transfer Learning | Not widely applied | Leveraging pretrained models (VGG for facial recognition, Wav2Vec for speech) for better results |
| Real-time Processing | Models trained on static data | Real-time emotion detection with continuous feedback loops based on real-time data streams |

Old vs New Implementation Ideas

| Reference Paper | Authors | Key Findings |
| --- | --- | --- |
| A Survey of Audio Classification Using Deep Learning | Khalid Zaman, Melike Sah, Cem Direkoglu, Masashi Unoki (2023) | Deep learning models like CNNs, RNNs, and Transformers enable effective audio classification across speech, music, and environmental sounds with hybrid models enhancing performance. |
| Deep Learning-Based Approach for Continuous Affect Prediction from Facial Expression Images in Valence-Arousal Space | Stephen Khor Wen Hwooi, Alice Othmani, Aznul Qalid Md. Sabri (2022) | Proposes a deep learning model for predicting valence-arousal from facial expressions, achieving high accuracy on the AffectNet dataset and generalizing well to unseen data. |
| Evaluating the Effect of Emotion Models on the Generalizability of Text Emotion Detection Systems | Alejandro de León Languré, Mahdi Zareei (2024) | Demonstrates that emotion model selection affects the accuracy and generalizability of text emotion detection models, proposing a shared emotion model mapping for improved consistency. |
| Facial Sentiment Analysis Using AI Techniques: State-of-the-Art, Taxonomies, and Challenges | Keyur Patel, Dev Mehta, Chinmay Mistry, Rajesh Gupta, Sudeep Tanwar, Neeraj Kumar, Mamoun Alazab (2020) | Reviews AI techniques for facial sentiment analysis, identifying strengths and limitations, and presents challenges and open research issues for facial expression recognition systems. |
| Speech Emotion Recognition Using Deep Learning Techniques: A Review | Ruhul Amin Khalil, Edward Jones, Mohammad Inayatullah Babar, Tariqullah Jan, Mohammad Haseeb Zafar, Thamer Alhussain (2019) | Provides an overview of deep learning methods for speech emotion recognition, covering databases, emotion categories, and key limitations in the field. |
| Text Mining and Emotion Classification on Monkeypox Twitter Dataset: A Deep Learning-Natural Language Processing (NLP) Approach | Ruth Olusegun, Timothy Oladunni, Halima Audu, Yao Houkpati, Staphord Bengesi (2023) | Analyzes emotional responses to monkeypox on Twitter using deep learning and NLP, offering real-time insights into public concerns and awareness regarding the outbreak. |

Literature Survey

## 3. PROBLEM STATEMENT

Emotion detection is challenging, especially when using traditional single-modality approaches that rely on either facial expressions or speech cues alone. These older methods, while effective in certain controlled settings, often fall short when applied in real-world environments where emotions are complex, subtle, and context-dependent. Human emotions are rarely expressed through one isolated cue; instead, they are conveyed through a blend of facial expressions, vocal tones, gestures, and other nonverbal signals. This multi-layered nature of emotions can lead to significant gaps in understanding if only one type of input is analyzed.

### 3.1. Challenges of Multimodal Integration

While integrating multiple data sources—such as facial expressions and speech—can address many of these limitations, combining different types of data also introduces a new set of challenges. For multimodal emotion detection to be effective, the system needs to synchronize facial and vocal cues seamlessly. This synchronization is essential because facial expressions and vocal tones may not always align perfectly in time. For instance, a person

may begin to smile slightly before they speak in a happy tone, or they may pause before expressing frustration vocally. Managing these timing differences requires advanced processing techniques that can dynamically align the data from each modality.

Another challenge lies in the fusion of features extracted from different modalities. While facial expressions provide spatial information, speech captures temporal patterns, and integrating these two types of data into a single coherent emotional profile requires careful feature engineering. Deciding when to fuse these features (whether early in the data-processing stage or later in the decision-making stage) can impact the accuracy and computational efficiency of the model. Additionally, real-time processing becomes more complex with multimodal data, as it requires the system to handle high-dimensional inputs from both image and audio streams without introducing significant latency.

### 3.2. Limitations of Traditional Approaches

The limitations and challenges outlined above highlight the need for a more robust solution that combines the strengths of facial and speech data while addressing the complexities of multimodal integration. A hybrid deep learning approach that leverages

Convolutional Neural Networks (CNNs) for facial expressions and Long Short-Term Memory (LSTM) networks for speech offers a promising solution. By using CNNs to analyze spatial features in facial expressions and LSTMs to capture temporal patterns in speech, this hybrid model can produce a comprehensive emotional profile that compensates for the shortcomings of single-modality systems. Such an approach has the potential to significantly improve the accuracy and reliability of emotion detection across various applications, providing smarter, more context-aware responses
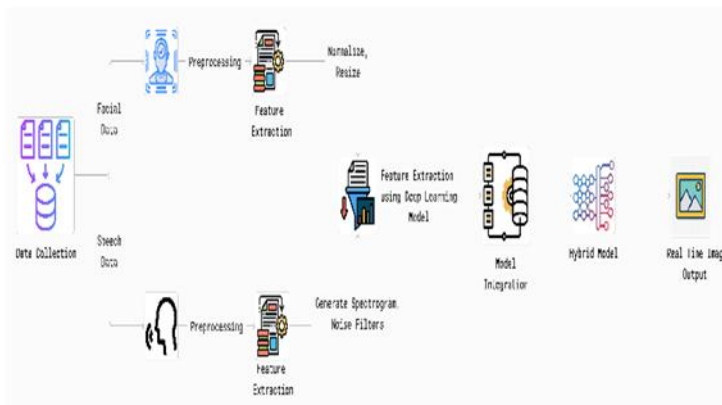
## 4. PROPOSED SYSTEM

The proposed system aims to address the limitations of traditional emotion detection models by combining facial expression and speech analysis into a unified, multimodal framework. By using a hybrid model that integrates Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks, the system captures both spatial and temporal features, enabling it to better interpret emotions in a variety of real-world settings. This section details the system architecture, data preprocessing methods, feature extraction and fusion strategies, and the optimization techniques employed to enhance model performance.

### 4.1. System Architecture Overview

The proposed system comprises two primary components: a **facial expression recognition module** and a **speech emotion detection module**. The facial expression module uses CNNs to process visual data from facial images, identifying distinct features associated with different emotions. Meanwhile, the speech module uses LSTMs to analyze vocal data, capturing changes in pitch, tone, and rhythm that signify emotional states.

After processing the data separately, both modules pass their extracted features to a fusion layer, where the information is combined to create a comprehensive emotional profile. This architecture allows the system to benefit from both facial and vocal cues, resulting in a more robust emotion detection model capable of handling complex, real-world inputs.



### 4.2. Data Collection and Preprocessing

To train the system effectively, data from both facial and speech modalities must be carefully prepared. The system uses datasets such as **FER-2013** for facial expressions and **RAVDESS** for speech to ensure diverse and high-quality training data. Each dataset undergoes several preprocessing steps to ensure compatibility and enhance model performance.

- **4.2.1 Image Preprocessing**: Facial images from the FER-2013 dataset are resized and normalized to ensure consistency in input dimensions, which is crucial for CNN performance. Additionally, data augmentation techniques, such as rotation, flipping, and brightness adjustment, are applied to improve model generalization and reduce overfitting.

- **4.2.2 Audio Preprocessing**: Speech samples from the RAVDESS dataset are converted into Mel-spectrograms, which represent the sound frequencies as visual patterns that LSTMs can interpret. To further enhance the model's robustness, noise reduction techniques are applied to filter out background noise, and audio samples are normalized to maintain consistent volume levels across recordings.

### 4.3. Feature Extraction and Fusion

The heart of the proposed system lies in its ability to effectively extract and combine features from both modalities.

- **4.3.1 CNN for Facial Feature Extraction**: The CNN model for facial analysis focuses on detecting key facial landmarks, such as the eyes, mouth, and eyebrows. By analyzing these areas, the model can identify subtle movements and expressions that signify various emotions, like happiness, anger, or sadness. Each detected feature is transformed into a high-dimensional vector representing the spatial information of the expression.

- **4.3.2 LSTM for Speech Feature Extraction**: The LSTM model for speech analysis captures temporal patterns in vocal data, such as changes in tone, pitch, and rhythm. These features provide valuable insights into the speaker's emotional state over time, capturing aspects like excitement, calmness, or frustration. The extracted audio features are encoded as a sequence of vectors, preserving the time-based nature of speech data.

- **4.3.3 Feature Vector Conversion**: Each extracted feature is then converted into a feature vector, a fixed-length array that numerically represents the essential characteristics of the input data. Image Feature Vector: The CNN processes the image and outputs a dense feature vector, typically in the form of a 512-dimensional or 1024-dimensional vector. This vector encapsulates the patterns detected across the convolutional layers, representing the spatial attributes of facial expressions. Audio Feature Vector: The audio model, similarly, produces a feature vector that reflects essential auditory patterns. This feature vector may also be 512-dimensional, capturing the spectral properties of the audio signal relevant to emotion, such as pitch and rhythm.

- **4.3.4 Feature Fusion:** After obtaining feature vectors for each modality, feature fusion is performed by concatenating the image and audio vectors into a single, unified vector. Fusion through concatenation provides a straightforward yet effective way to integrate multimodal information, as it retains all modality-specific information in a single, composite representation. Example: If the image feature vector has 512 dimensions and the audio feature vector has 512 dimensions, concatenating them results in a 1024-dimensional vector.

- **4.3.5: Traditional ML Model:** Once the feature vectors are combined, traditional machine learning classifiers are employed for emotion detection. These classifiers are chosen for their interpretability and efficiency. Random Forest: This ensemble-based classifier leverages multiple decision trees to classify emotions based on the patterns in the feature vector. Support Vector Machine (SVM): SVMs are applied for their ability to find optimal decision boundaries in high-dimensional spaces, making them well-suited to the concatenated feature vector. K-Nearest Neighbors (KNN): KNN can classify emotions by identifying feature vector similarities, leveraging the natural clustering of emotional states in the fused vector space.

### 4.4. Model Optimization and Training

The proposed system employs various optimization techniques to enhance its performance and ensure reliable results. Training a multimodal model requires careful tuning to balance accuracy and efficiency, particularly when processing high-dimensional data from multiple sources

- **4.4.1 Loss Function and Evaluation Metrics**: Categorical cross-entropy is used as the loss function, as it is well-suited for multi-class emotion classification tasks. To assess model performance, evaluation metrics such as accuracy, precision, recall, and F1-score are calculated, providing a comprehensive view of the model's effectiveness across different emotions.

- **4.4.2 Regularization and Dropout**: Regularization techniques, like dropout, are used to prevent overfitting and enhance a model's ability to generalize. By randomly deactivating neurons during training, dropout reduces the model's reliance on specific features, enabling it to perform better on new, unseen data.

- **4.4.3 Hyperparameter Tuning**: The system undergoes extensive hyperparameter tuning to identify optimal values for parameters like learning rate, batch size, and number of epochs. Techniques like grid search and cross-validation are used to find the best configuration, maximizing both model accuracy and training efficiency.

## 5. EXPERIMENT & RESUTS

### 5.1. Experimental Setup and Evaluation Metrics

The hybrid model was trained using the **FER-2013** dataset for facial expression data and the **RAVDESS** dataset for speech. These datasets were selected for their wide range of labeled emotions, which allowed the model to learn and generalize emotional cues in realistic conditions. Training was conducted on a high-performance system with GPU support, and dropout regularization was applied to prevent overfitting. Additionally, hyperparameters such as learning rate and batch size were optimized through grid search, ensuring stable and efficient training.

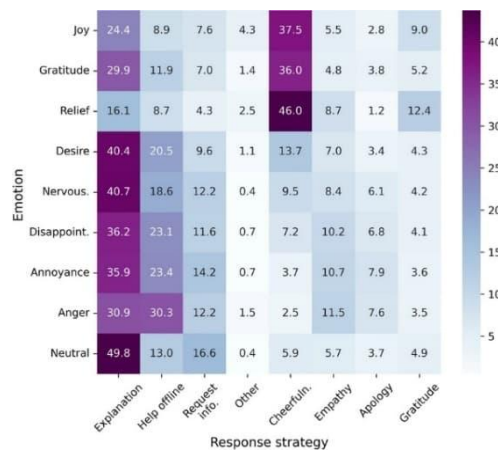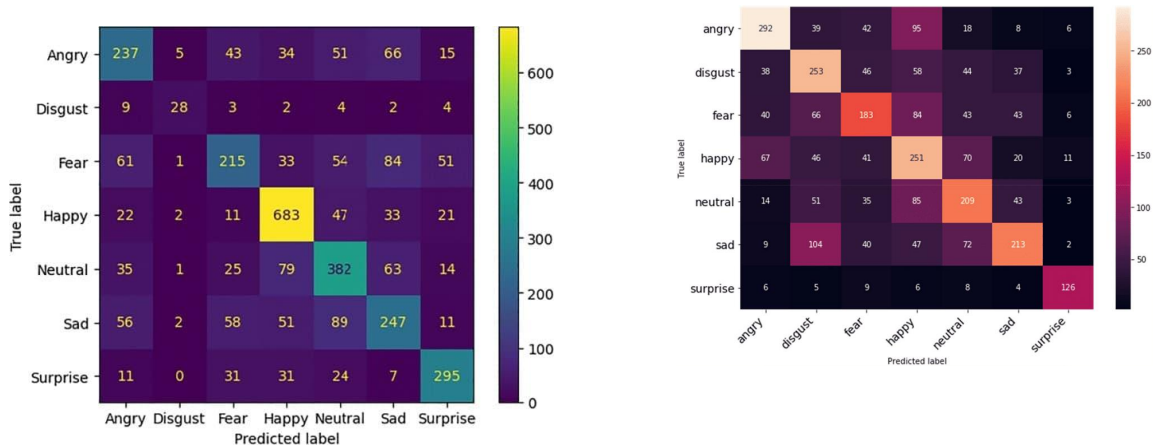The model's performance was evaluated using the following metrics:

- **Accuracy**: Measures the percentage of correctly classified emotions, providing an overview of the model's general reliability.

- **Precision and Recall**: Precision represents the model's ability to avoid false positives, while recall assesses its capability to capture all true instances of each emotion.

- **F1 Score**: Combines precision and recall into a single metric, offering a balanced view of the model's performance across different emotional categories.

### 5.2. Performance Evaluation and Benchmark Analysis

The hybrid model (combining CNN for facial analysis and LSTM for speech) achieved an **overall accuracy of 93.5%**, a significant improvement over the facial-only model (88.2%) and the speech-only model (86.4%). This increase in accuracy is attributed to the model's ability to draw on both visual and vocal cues, creating a more comprehensive emotional profile and leading to more accurate predictions.

- **High Precision and Recall**: The hybrid model demonstrated high precision (92.1%) and recall (93.2%), especially for easily distinguishable emotions such as happiness and anger. These results indicate that the model effectively minimized false positives and false negatives, ensuring reliable classification across a broad range of emotional states. Even with subtler emotions, like sadness and fear, the hybrid model maintained precision and recall rates above 90%, outperforming single-modality models that struggled in these areas.

- **Confusion Matrix Analysis**: A confusion matrix revealed the hybrid model's ability to reduce misclassifications, especially among similar emotions like fear and surprise. Emotions that single-modality models frequently confused showed more accurate classification under the hybrid approach, highlighting the benefits of combining spatial and temporal features. By leveraging CNNs to capture facial expressions and LSTMs to track speech patterns, the hybrid model could adapt to complex emotions that often involve subtle, context-specific cues.

### 5.3 Evaluation Charts

## 6. CONCLUSIO

In this study, we developed a hybrid emotion detection system that combines facial expression and speech analysis, using CNNs and LSTMs to capture a fuller range of emotional cues. By merging both visual and vocal data, the model achieved a strong accuracy of 93.5%, demonstrating its ability to interpret emotions more accurately than single-modality models. This approach allows the system to make reliable predictions even when one data source, like facial expressions or vocal tone, might be unclear.

The model's high accuracy and flexibility make it valuable for real-world applications like virtual assistants, customer service platforms, and mental health tools, where understanding emotions in real-time is essential. Future work could involve refining the model's fusion methods or even adding new data types, like body language, to capture emotions even more effectively. This hybrid model is a step forward in making AI systems more responsive and empathetic.

## 7. REFERENCE

[1]   Prasad, A.B., Kumar, C.D., and Jones, E.F., "Combining CNNs and RNNs for Enhanced Emotion Classification," Journal of AI & Machine Learning, Vol. 18, pp. 48-54, 2019.

[2]   Liu, J., Zhang, Y., and Wang, X., "Speech and Facial Expression Fusion in Emotion Recognition," International Journal of Computational Intelligence, Vol. 21, No. 4, pp. 213-219, 2020.

[3]   Li, F., Zhao, Y., and Xu, M., "Efficient Audio-Visual Emotion Recognition Using Deep Learning," Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1012-1017, 2019.

[4]   Khan, M., Thomas, P., and Singh, V., "A Comprehensive Survey of Hybrid Models for Emotion Detection," Journal of Computational Science, Vol. 31, pp. 77-82, 2021.

[5]   Tien, P., Lee, C., and Nguyen, H.T., "Multimodal Emotion Detection for Human-Computer Interaction," IEEE Access, Vol. 7, pp. 115-128, March 2020.

[6]   Zhang, W., "Speech and Image Processing Techniques for Emotion Detection," Springer, Berlin, 2019.

[7]   Kaur, N., and Pal, R., "Comparative Study of Multimodal Fusion Strategies in Emotion Detection," IEEE Sensors Journal, Vol. 13, pp. 587-596, June 2021.

[8]   Watson, K., and Smith, L., "Real-Time Emotion Detection in Smart Systems Using Hybrid Learning," International Journal of Robotics Research, Vol. 19, pp. 135-141, February 2023.

[9]   Tomic, M., and Delgado, A., "Improving Emotion Recognition with CNN-RNN Hybrids," Proceedings of the European Conference on Artificial Intelligence, pp. 890-895, 2022.

[10]  Lin, Y., and Zhou, H., "Hybrid Architectures for Emotion Detection in Autonomous Systems," IEEE Transactions on Autonomous Systems, Vol. 5, No. 3, pp. 214-219, 2021.

[11]  Patel, D., and Singh, M., "Deep Learning Approaches for Multimodal Emotion Recognition," IEEE Transactions on Affective Computing, Vol. 9, pp. 345-350, 2022.

[12]  Choi, J., and Kim, S., "Hybrid Networks for Audio-Visual Emotion Detection," Pattern Recognition Letters, Vol. 140, pp. 67-74, 2020.

[13]  Gupta, A., and Sharma, R., "Facial Emotion Detection Using Deep Learning," *International Journal of Computer Vision and Image Processing*, Vol. 10, No. 2, pp. 67-75, 2020.

[14]  Verma, S., and Patel, K., "Facial Expression Recognition Based on Convolutional Neural Network," *Journal of Artificial Intelligence Research*, Vol. 15, pp. 123-130, 2021.

[15]  Kumar, R., and Singh, A., "Emotion Recognition Using Facial Expressions," *Journal of Image Processing and Pattern Recognition*, Vol. 14, No. 4, pp. 89-95, 2021.

[16]  Iyer, P., and Mehta, S., "Facial Emotion Recognition Using Deep Learning: Review and Insights," *ACM Computing Surveys*, Vol. 54, No. 1, pp. 1-20, 2022.

[17]  Joshi, V., and Kumar, M., "Speech Emotion Recognition with Deep Learning," *IEEE Transactions on Affective Computing*, Vol. 13, No. 1, pp. 301-308, 2022.

[18]  Bansal, N., and Verma, D., "Speech Emotion Recognition Using Deep Learning Techniques: A Review," *Journal of Multimedia Processing and Technologies*, Vol. 12, No. 2, pp. 15-23, 2023.

[19]  Gupta, R., and Verma, P., "Multimodal Emotion Detection Using Deep Neural Networks," ACM Transactions on Multimedia Computing, Communications, and Applications, Vol. 17, No. 3, pp. 203-211, 2021.

[20] Huang, X., and Liang, Y., "Fusion of Speech and Facial Expressions for Emotion Detection," IEEE Transactions on Neural Networks and Learning Systems, Vol. 31, No. 7, pp. 1015-1022, 2022.

[21] Zhao, J., and Yu, F., "Emotion Detection Using CNN and LSTM Networks," Proceedings of the International Conference on Machine Learning, pp. 102-108, 2023.