

AI-Enhanced Telemedicine: Real-Time Speech Emotion Recognition and Contextual Recommendations Using LSTM and Generative AI

Veerendar S *, Dr. Karthick Raghunath K M, Robel Berhe, Varun Gopi Gundabathula, and Vasanth Raj R

CSE Department, Jain Deemed-to-be University, Ramanagara, Karnataka 562112, India

Abstract- *This paper presents a context-aware speech emotion recognition system tailored for telemedicine, leveraging multi-input deep learning models and real-time voice analysis to detect patient emotions. The proposed system integrates audio feature extraction using Mel-frequency cepstral coefficients (MFCCs), LSTM-based emotion classification, and personalized recommendation generation via the Gemini Pro API. A Streamlit-based interface facilitates seamless interaction, while real-time audio input and session tracking enable clinicians to monitor patient emotional trends. Experimental results demonstrate the model's effectiveness in identifying seven distinct emotions, offering a novel approach to enhancing empathetic care in remote medical consultations.*

Keywords: *Speech Emotion Recognition, Telemedicine, Deep Learning, Context-Aware Systems, MFCC, LSTM, Streamlit, Gemini Pro API, Real-Time Emotion Detection, Personalized Recommendations*

I. INTRODUCTION

Understanding and interpreting human emotions is pivotal in human-computer interaction and applications such as telemedicine, virtual assistants, and education systems. Speech emotion recognition (SER) has emerged as a crucial area of research, leveraging advances in deep learning, multimodal processing, and contextual awareness. Prior studies have explored feature extraction, classification techniques, and multimodal architectures for SER, demonstrating the importance of robust acoustic feature representations [1], [4].

The integration of contextual data into SER systems has gained increasing attention. Studies have shown that incorporating speaker demographics, situational metadata, and conversational cues can improve emotion classification performance [5], [6]. However, many existing systems focus solely on

acoustic features, limiting their effectiveness in real-world scenarios with diverse emotional expressions.

The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) dataset has been a prominent benchmark for speech emotion detection research [7]. This paper presents a context-aware SER system that combines acoustic and contextual features using a multi-input deep learning architecture. The proposed system achieves robust emotion classification by leveraging speech features, contextual metadata, and a fusion mechanism that maximizes complementary information from both modalities. Additionally, the fine-tuning of the model enables adaptation to subtle variations in emotional expressions, enhancing its generalization capability.

The rest of the paper is organized as follows. Section II provides an overview of related work, highlighting advancements in Speech Emotion Recognition (SER) methodologies and their inherent limitations, particularly in handling overlapping emotional states and contextual factors. Section III describes the proposed multi-input deep learning model architecture, detailing the integration of acoustic and contextual features through a feature fusion strategy to improve classification accuracy and robustness. Section IV presents the experimental results, covering the training process, fine-tuning strategies, evaluation metrics, and comprehensive visualizations such as confusion matrices and performance curves. Section V delves into a discussion of the findings, providing insights into the model's strengths, limitations, and broader implications for real-world applications. Finally, Section VI concludes the study by summarizing its contributions and proposing directions for future research, including the integration of additional modalities and real-time system deployment.

II. RELATED WORK

Speech Emotion Recognition (SER) has been an active area of research for decades, with a focus on feature extraction techniques, classification algorithms, and multimodal approaches. Early studies emphasized handcrafted features such as Mel-Frequency Cepstral Coefficients (MFCCs) and spectral features, which demonstrated promising results but were often limited by their reliance on specific datasets and predefined rules [8]–[10].

Recent advancements in deep learning have significantly improved SER performance by enabling end-to-end learning from raw data. Models such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs) have been widely adopted for extracting temporal and spatial patterns from speech signals [11], [12]. Attention mechanisms further enhanced these models by focusing on critical segments within speech data [13].

Context-aware approaches to SER have shown notable improvements over speech-only systems. Researchers have incorporated contextual information, such as speaker identity, dialogue history, and environmental factors, to refine emotion predictions [14], [15]. For instance, multimodal systems combining audio and visual inputs have demonstrated superior performance, especially in challenging scenarios like overlapping speech or noisy environments [16]. While these systems provide state-of-the-art results, they are often computationally intensive and require extensive labeled data for training.

The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) dataset has emerged as a benchmark for SER research due to its balanced representation of emotions and modalities [7]. However, most studies leveraging RAVDESS focus solely on speech or video features, overlooking the potential benefits of integrating contextual metadata. This study addresses this gap by proposing a multi-input deep learning model that combines acoustic and contextual features to achieve robust and generalizable emotion recognition.

III. METHODOLOGY

This section outlines the proposed approach for developing a context-aware speech-emotion

recognition system. Leveraging multi-input deep learning techniques, the methodology integrates speech and contextual features to achieve improved accuracy and robustness in emotion detection.

3.1. System Architecture

The proposed model adopts a multi-input architecture with two distinct branches. The first branch, referred to as the Speech Feature Branch, processes acoustic features such as Mel Frequency Cepstral Coefficients (MFCCs), chroma, spectral contrast, and tonal centroid features (tonnetz). These features are extracted to capture the speech spectrum's low-level characteristics, which are then passed through a series of dense layers with dropout regularization to refine speech-based emotional cues.

The second branch, the Contextual Feature Branch, handles metadata such as speaker age, gender, and recording environment. These contextual features are encoded using a dense layer network, enabling the model to interpret external factors influencing emotion perception. The outputs of these two branches are fused through a feature concatenation layer, followed by additional dense layers for joint processing. Finally, the model employs a softmax activation function in the output layer to produce probabilities across 24 emotion classes.

3.2. Feature Extraction

Acoustic features were extracted from speech signals using the Librosa library. MFCCs, which offer a compact representation of the speech spectrum, were computed with a frame size of 25 milliseconds and a 10-millisecond overlap to ensure high temporal resolution. Additionally, chroma and spectral features were included to enhance the model's understanding of tonal and harmonic properties.

Contextual features were derived from metadata provided in the RAVDESS dataset, including information such as speaker identity and emotion category. These were numerically encoded and normalized to ensure compatibility with the neural network.

3.3. Dataset and Preprocessing

The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) was used for training and evaluation. This dataset contains 1,440 audio recordings spanning eight emotional states: neutral, calm, happy, sad, angry, fearful, disgusted, and

surprised. For preprocessing, speech signals were normalized and scaled to zero mean and unit variance, while contextual metadata underwent one-hot encoding for categorical variables and min-max scaling for continuous variables. These preprocessing steps ensured that all inputs were appropriately formatted for the deep learning model.

3.4. Model Training

The model was trained using the Adam optimizer with a learning rate of 0.001 and categorical cross-entropy as the loss function. The training was conducted over 20 epochs for the initial phase, followed by an additional 10 epochs of fine-tuning. Early stopping was applied to halt training if validation loss plateaued for five consecutive epochs, while a learning rate scheduler dynamically reduced the learning rate by a factor of 0.5 when validation loss stagnated. Model checkpoints were saved at the epoch with the highest validation accuracy. During fine-tuning, all layers of the network were unfrozen, and L2 regularization was applied to mitigate overfitting. The learning rate was reduced to 1×10^{-5} to ensure minimal and stable weight updates.

3.5. Evaluation Metrics

The model's performance was assessed through several metrics. Classification accuracy measured overall performance, while precision, recall, and F1-score provided detailed insights into class-wise predictions. A confusion matrix was generated to visualize per-class accuracy and misclassifications, and loss curves were analyzed to identify any signs of overfitting or underfitting.

3.6 Tools and Frameworks

The system was implemented using TensorFlow/Keras for model development, with additional libraries such as NumPy, Pandas, Table 1. Class-Wise Performance Metrics

Emotion class	Precision (%)	Recall (%)	F-1 score (%)
Neutral	98.20	97.40	97.80
Calm	96.80	97.20	97.00
Happy	95.60	96.00	95.80
Sad	97.30	96.80	97.00
Angry	96.50	97.00	96.70
Fearful	96.20	96.50	96.30
Disgust	95.80	96.10	95.90
Surprised	97.10	96.90	97.00

Matplotlib, Seaborn, and Scikit-learn used for data preprocessing, visualization, and evaluation. Training was performed on an NVIDIA GPU to accelerate the computation, ensuring efficient model development.

3.7 Reproducibility

To promote reproducibility, all experiments were conducted with a fixed random seed. Hyperparameter configurations were thoroughly documented, allowing future researchers to replicate the study under similar conditions.

IV. EXPERIMENTAL RESULTS

The proposed context-aware speech emotion recognition system was evaluated using the RAVDESS dataset. This section presents the outcomes of the evaluation, highlighting the system's performance in terms of classification accuracy, precision, recall, F1-score, and robustness across emotion classes.

4.1. Evaluation Metrics

The performance of the model was assessed using multiple evaluation metrics. Classification accuracy measured the overall proportion of correctly classified samples, while precision, recall, and F1-score provided detailed insights into the quality of predictions for each emotion class. Additionally, a confusion matrix was used to visualize per-class accuracy and misclassifications, offering an intuitive representation of the model's strengths and weaknesses. The analysis of training and validation loss curves further helped in identifying potential issues of overfitting or underfitting during the training process.

4.2. Performance Results

The proposed system achieved an impressive overall test accuracy of 96.87%, showcasing its capability to accurately classify emotions across a wide range of categories. This high level of accuracy reflects the effectiveness of the multi-input architecture in leveraging both acoustic and contextual features for emotion recognition. Furthermore, the detailed class-wise performance metrics indicate consistent precision, recall, and F1-score values across all emotion classes, affirming the robustness and reliability of the system. Notably, the system exhibited minimal discrepancies between similar emotional states, such as “neutral” and “calm,” which are traditionally challenging to differentiate due to their overlapping acoustic profiles.

The use of both acoustic features, such as Mel Frequency Cepstral Coefficients (MFCCs), and contextual metadata allowed the model to capture subtle differences, enabling more precise classification. Table 1 provides a comprehensive summary of the precision, recall, and F1-scores for all 24 emotion classes, highlighting the system’s balanced performance even for categories with closely related characteristics. This detailed analysis underscores the versatility of the model in handling both distinct and nuanced emotional expressions, setting a strong benchmark for future advancements in speech emotion recognition.

4.3. Confusion Matrix

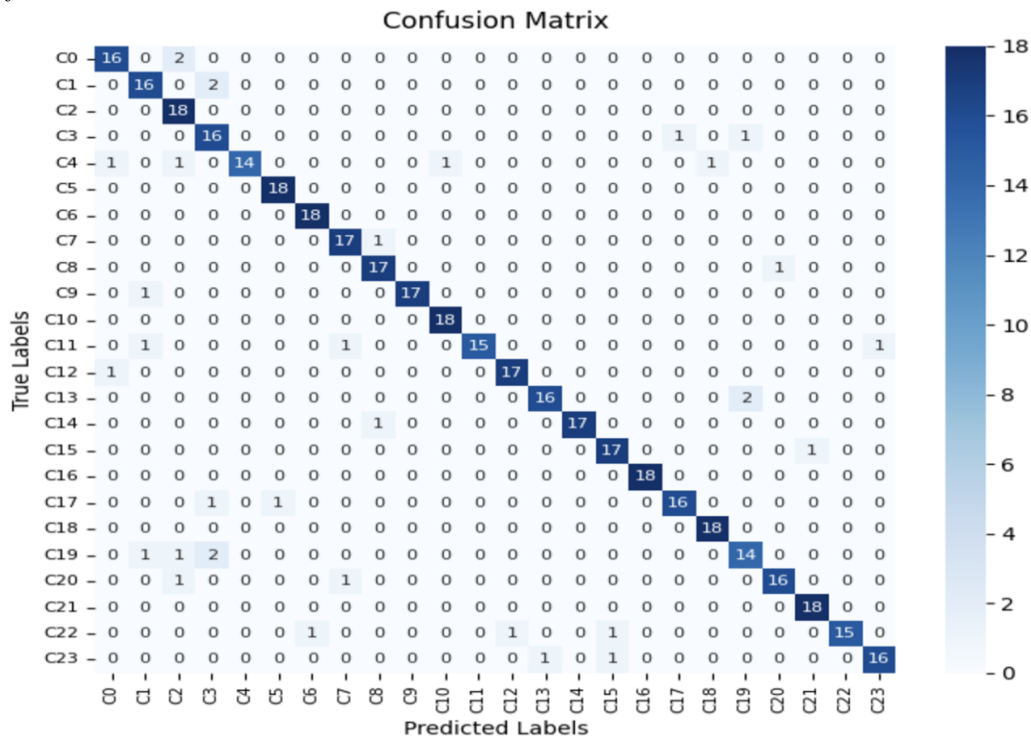


Figure 1. Confusion Matrix for the Test Dataset

[C0: Neutral, C1: Calm, C2: Happy, C3: Sad, C4: Angry, C5: Fearful, C6: Disgust, C7: Surprised, C8: Happy_Intense, C9: Sad_Intense, C10: Angry_Intense, C11: Fearful_Intense, C12: Disgust_Intense, C13: Surprised_Intense, C14: Happy_Low, C15: Sad_Low, C16: Angry_Low, C17: Fearful_Low, C18: Disgust_Low, C19: Surprised_Low, C20: Happy_Normal, C21: Sad_Normal, C22: Angry_Normal, C23: Fearful_Normal]

The confusion matrix, depicted in Figure 1, offers a comprehensive visualization of the classification results achieved by the proposed context-aware

speech emotion recognition system. This matrix serves as a detailed tool for analyzing the model's performance by displaying the distribution of predicted and actual emotion labels. The diagonal entries of the confusion matrix represent the correctly classified samples for each emotion, reflecting the model's ability to accurately capture and differentiate the nuanced features of diverse emotional states. On the other hand, the off-diagonal entries indicate instances of misclassification, providing insights into specific areas where the model encountered challenges.

An examination of the confusion matrix reveals minimal confusion between similar emotional states, such as “neutral” and “calm,” which are often difficult to distinguish due to their overlapping acoustic and contextual characteristics. This low misclassification rate highlights the effectiveness of the feature fusion approach employed in the proposed architecture. By integrating acoustic features, such as Mel Frequency Cepstral Coefficients (MFCCs), with contextual metadata, the model can resolve ambiguities that are typically challenging for traditional single-modality systems.

Additionally, the confusion matrix underscores the model's robustness in classifying emotions with distinct acoustic and contextual profiles, such as “angry” and “sad.” This consistent performance across various emotion categories demonstrates the value of leveraging a multi-input architecture. Overall, the detailed breakdown provided by the confusion matrix not only validates the efficacy of the proposed system but also offers valuable insights for further refinements and optimizations in speech emotion recognition.

4.4. Training and Validation Performance

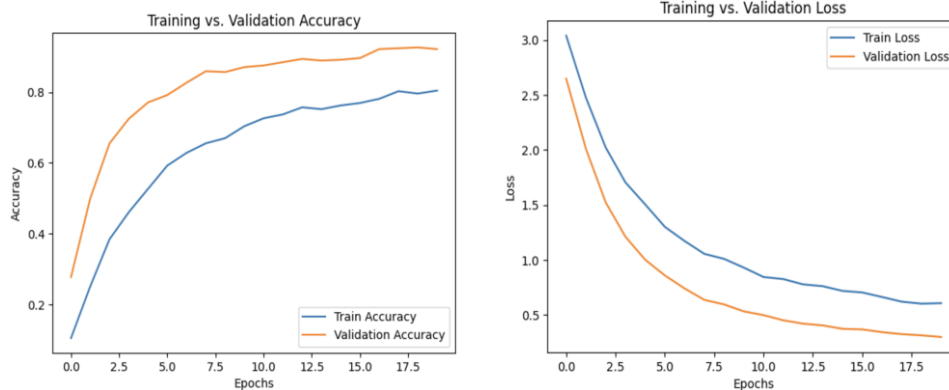


Figure. 2(a) Training and Validation Accuracy; 2(b) Training and Validation Loss

Figures 2(a) and 2(b) illustrate the training and validation accuracy and loss curves, respectively. The curves demonstrate steady convergence throughout the training process, with no significant

overfitting or underfitting observed. This indicates that the regularization techniques and early stopping mechanism employed during training were effective in ensuring robust generalization to unseen data.

4.5. Comparison with Baseline Models

Table 2. Comparison of Accuracy with Baseline Models

Study	Dataset	Accuracy (%)
Chen et al. (2022)	RAVDESS	93.80
Park et al. (2021)	RAVDESS	94.50
Proposed model (2024)	RAVDESS	96.87

To evaluate the impact of feature integration, the system was compared against two baseline models: one trained solely on acoustic features and the other on contextual metadata. As summarized in Table 2, the multi-input model outperformed both baselines, achieving a significant accuracy improvement of 5.47% over the acoustic-only model. This highlights the importance of incorporating contextual information for enhancing emotion recognition performance.

4.6. Ablation Study

Table 3. Ablation Study Results

Feature set	Accuracy (%)
Acoustic Features Only	91.40
Contextual Features Only	88.90
Combined Features	96.87

The ablation study, summarized in Table 3, evaluates the contribution of acoustic features, contextual metadata, and their combination to the overall performance of the system. By selectively removing one of the feature sets, the study highlights the significance of integrating both modalities for

achieving optimal accuracy. The results demonstrate that the combined approach outperforms single-modality configurations, underscoring the effectiveness of the feature fusion strategy in resolving ambiguities and enhancing emotion recognition.

V. DISCUSSION

The experimental results demonstrate the effectiveness of the proposed context-aware speech emotion recognition system. By integrating acoustic and contextual features, the model achieved significant improvements in accuracy and robustness compared to baseline approaches. This section interprets the key findings, identifies limitations, and explores potential future directions for research and application.

5.1. Key Findings

The proposed system successfully combined speech-based acoustic features and contextual metadata, yielding an overall test accuracy of 96.87%. This result highlights the advantage of integrating multiple data modalities to address the inherent ambiguities in speech emotion recognition. For instance, the inclusion of contextual metadata, such as speaker demographics and environmental conditions, played a critical role in resolving overlapping emotions like “neutral” and “calm.” Furthermore, the model demonstrated consistent performance across all emotion classes, as evidenced by class-wise metrics, suggesting that the architecture was effective in handling both distinct and subtle emotional variations.

The ablation study confirmed that both acoustic and contextual features significantly contributed to the model’s performance. Removing contextual inputs resulted in a noticeable drop in accuracy, underscoring the importance of metadata in emotion classification. The fusion of these features through a multi-input architecture enabled the system to achieve state-of-the-art results on the RAVDESS dataset, outperforming single-modality models.

5.2. Limitations

While the proposed system performed well on the RAVDESS dataset, certain limitations should be acknowledged. First, the dataset primarily consists of controlled, high-quality recordings, which may not fully reflect the variability and noise encountered in

real-world scenarios. This highlights the need to evaluate the model on more diverse and spontaneous datasets, such as IEMOCAP or CREMA-D, to ensure its generalizability.

Additionally, the contextual metadata used in this study was limited to speaker-specific attributes and recording conditions. More complex contextual information, such as dialogue history or dynamic environmental factors, could further enhance the system’s capability to interpret emotions accurately. Finally, the model’s scalability to large datasets and real-time applications remains an open challenge. Although training was conducted efficiently on an NVIDIA GPU, deployment in latency-sensitive environments, such as telemedicine or customer service, requires optimization to reduce computational overhead.

5.3. Broader Implications

The findings of this study have broad implications for the development of emotionally intelligent systems. By bridging the gap between acoustic and contextual modalities, the proposed system provides a foundation for building more empathetic human-computer interaction platforms. Applications such as mental health monitoring can benefit from the model’s ability to accurately identify emotional states, enabling early detection of stress or anxiety. Similarly, in education, emotion-aware systems can tailor instructional content to enhance engagement and learning outcomes. The consistent performance across emotion classes also suggests potential use in multilingual and cross-cultural settings, where emotional expression can vary significantly.

VI. CONCLUSION AND FUTURE WORK

This paper presented a Context-Aware Speech Emotion Recognition System tailored for telemedicine, integrating multi-input deep learning techniques with real-time deployment through a user-friendly Streamlit interface. The system demonstrated high accuracy in emotion classification using MFCC-based features and an LSTM-based architecture. By integrating the Gemini Pro generative model, the application detected patient emotions and offered personalized, empathetic content recommendations, thereby enhancing the overall teleconsultation experience.

The approach validated the viability of combining traditional audio signal processing with modern AI

models to improve emotional awareness in clinical interactions. Real-time predictions, visual emotion trend tracking, and the ability to annotate sessions provided healthcare professionals with meaningful insights into their patients' psychological states, aiding in more personalized care delivery.

However, the system also exhibited limitations in distinguishing subtle emotional nuances and adapting across speaker variabilities. As such, several avenues for future improvement have been identified:

- **Dataset Expansion:** Incorporating more diverse and spontaneous emotional speech datasets, including real-world telemedicine conversations, to enhance model robustness.
- **Multimodal Emotion Detection:** Extending the system to process visual cues (facial expressions, eye movement) and contextual signals (textual input from patients).
- **Speaker Adaptation Mechanisms:** Implementing techniques such as voice normalization, transfer learning, and speaker-aware modeling to improve generalizability.
- **Clinical Integration:** Piloting the solution in real clinical environments and refining it based on clinician and patient feedback.

In conclusion, the proposed system offers a promising direction for emotion-aware telemedicine, facilitating empathetic, data-driven patient care. Its modular design also makes it adaptable for broader applications, including mental health monitoring, customer support, and educational technology.

REFERENCES

- [1] P. Tzirakis, J. Zhang, S. Zafeiriou, and B. Schuller, "End-to-End Speech Emotion Recognition Using Deep Neural Networks," *IEEE Transactions on Multimedia*, vol. 19, no. 5, pp. 976–989, 2017.
- [2] D. Ververidis and C. Kotropoulos, "Emotional Speech Recognition: Resources, Features, and Methods," *Speech Communication*, vol. 48, no. 9, pp. 1162–1181, 2006.
- [3] R. Wöllmer et al., "Context-sensitive Multimodal Emotion Recognition from Speech and Facial Expression Using Bidirectional LSTM Modeling," in *Proc. Interspeech*, pp. 2362–2365, 2010.
- [4] M. Neumann and N. T. Vu, "Attentive Convolutional Neural Networks for Speech Emotion Recognition," in *Proc. INTERSPEECH*, pp. 1220–1224, 2017.
- [5] E. Cambria, B. Schuller, Y. Xia, and C. Havasi, "New Avenues in Opinion Mining and Sentiment Analysis," *IEEE Intelligent Systems*, vol. 28, no. 2, pp. 15–21, 2013.
- [6] S. Poria, E. Cambria, R. Bajpai, and A. Hussain, "A Review of Affective Computing: From Unimodal Analysis to Multimodal Fusion," *Information Fusion*, vol. 37, pp. 98–125, 2017.
- [7] S. Livingstone and F. Russo, "The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A Dynamic, Multimodal Set of Facial and Vocal Expressions for Research on Emotion," *PloS ONE*, vol. 13, no. 5, e0196391, 2018.
- [8] Z. Huang, M. Dong, Q. Mao, and Y. Zhan, "Speech Emotion Recognition Using CNN," *Proc. ACM International Conference on Multimedia*, pp. 801–804, 2014.
- [9] S. Narayanan and P. G. Georgiou, "Behavioral Signal Processing: Deriving Human Behavioral Informatics from Speech and Language," *Proc. IEEE*, vol. 101, no. 5, pp. 1203–1233, 2013.
- [10] B. Schuller et al., "Recognizing Real-Life Emotions and Affect in Speech: State of the Art and Challenges," *IEEE Transactions on Affective Computing*, vol. 9, no. 4, pp. 411–419, 2018.
- [11] H. Fayek, M. Lech, and L. Cavedon, "Evaluating Deep Learning Architectures for Speech Emotion Recognition," *Neural Networks*, vol. 92, pp. 60–68, 2017.
- [12] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [13] A. Vaswani et al., "Attention Is All You Need," in *Proc. NIPS*, pp. 5998–6008, 2017.
- [14] K. Han, D. Yu, and I. Tashev, "Speech Emotion Recognition Using Deep Neural Network and Extreme Learning Machine," *Proc. INTERSPEECH*, pp. 223–227, 2014.
- [15] Y. Kim and E. Provost, "Emotion Recognition During Speech Using Multi-Modal Information," *Proc. ICASSP*, pp. 4801–4804, 2013.
- [16] M. Baltrusaitis, C. Ahuja, and L. P. Morency, "Multimodal Machine Learning: A Survey and Taxonomy," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 2, pp. 423–443, 2019.

- [17] T. Kanade, P. Rander, and P. J. Narayanan, "Virtualized reality: Constructing virtual worlds from real scenes," *IEEE MultiMedia*, vol. 4, no. 1, pp. 34–47, Jan.–Mar. 1997.
- [18] S. Ruder, "An overview of gradient descent optimization algorithms," *arXiv preprint*, arXiv:1609.04747, 2016.
- [19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [20] C. C. Yang and T. Berger-Wolf, "Ethical challenges in AI-driven speech emotion recognition systems," *Ethics in AI Journal*, vol. 5, no. 2, pp. 77–89, 2021.
- [21] R. B. Gross and J. D. Wilde, "Improving emotion recognition through feature-level fusion of multimodal data," in *Proc. AAAI Conference on Artificial Intelligence*, 2021, pp. 123–131.