

Personalized Emotion Detection Adapting Models to Individual Emotional Expressions

¹Diwakar Mainali; ²Saraswoti Shrestha; ³Umesh Thapa; ⁴Sanjib Nepali
Tribhuvan University, Humanities and Social science, BCA

Abstract:- Emotion recognition from text and speech has become a critical area of research in artificial intelligence (AI), enhancing human-computer interaction across various sectors. This paper explores the methodologies used in emotion recognition, focusing on Natural Language Processing (NLP) for text and acoustic analysis for speech. It reviews key machine learning and deep learning models, including Support Vector Machines (SVM), neural networks, and transformers, and highlights the datasets commonly used in emotion detection studies. The paper also addresses challenges such as multimodal integration, data ambiguity, and ethical considerations like privacy concerns and bias in models. Applications in customer service, healthcare, education, and entertainment are discussed, showcasing the growing importance of emotion recognition in AI-driven systems. Future research directions, including advancements in deep learning, multimodal systems, and real-time processing, are also explored to address existing limitations.

Keywords:- Emotion Recognition, Natural Language Processing (NLP), Acoustic Analysis, Machine Learning, Deep Learning, Speech Emotion Recognition, Text Emotion

Recognition, Multimodal Systems, AI Applications, Ethical Considerations.

I. INTRODUCTION

Emotion recognition is recognising human emotions via facial expressions, body language, voice, and writing. It aims to bridge computer-human interactions as part of AI and ML. Effective communication requires a strong understanding of emotions since they say more than words. Emotion detection improves user experience in healthcare and customer service because emotional cues affect decision-making, social interactions, and behaviour [1]. Spoken and written language emotion recognition is unique and important. Analysing word choice, phrase structure, and punctuation can assist extract emotions from text. Sentiment analysis and NLP are widely used to interpret these signals and classify emotions like joy, sadness, rage, and surprise. Similarly, speech emotion identification emphasises voice features like intensity, rhythm, pitch, and tone. Even without language, speech acoustics may convey a lot of emotion [2]. Machine learning algorithms analyse vocal features to determine emotions. Speech emotion recognition is used in call centres to improve customer service by understanding clients' emotions.

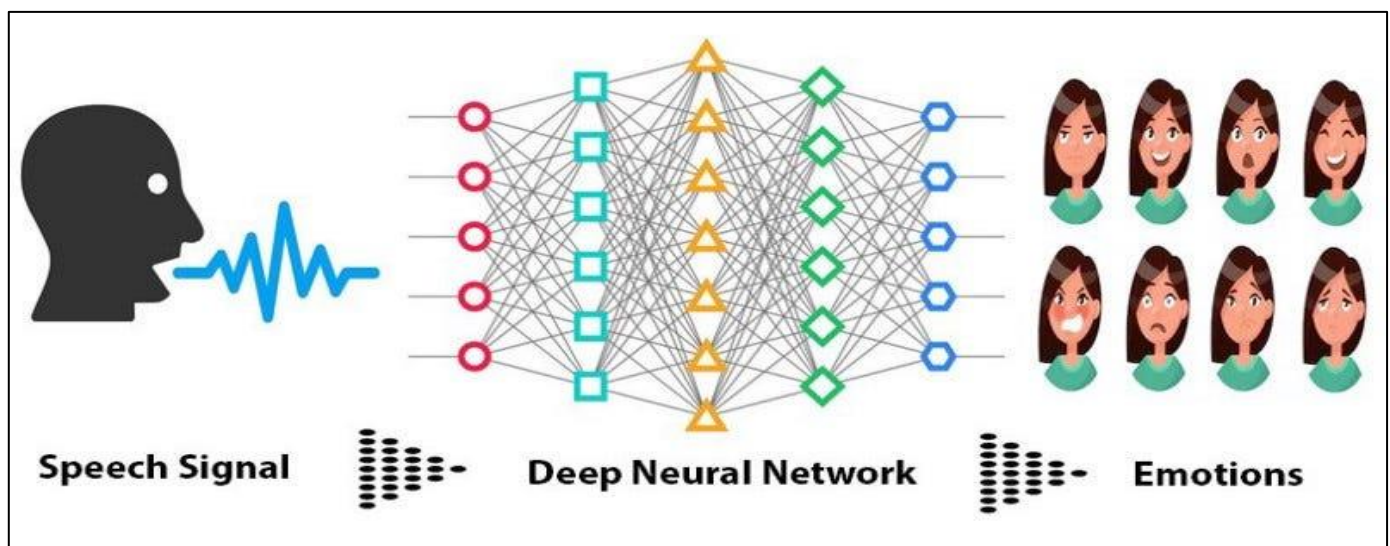


Fig 1: Emotion Recognition

This study reviews text and spoken emotion recognition literature to synthesise past work. This lecture will discuss emotion detection methods from both domains, the challenges of accurately identifying emotions, and its real-world applications. The study analyses secondary sources to

assess current research. It compares methodologies, examines improvements, and suggests future growth. Improving human-computer interactions demands keeping up with emotion identification advances.

II. BACKGROUND AND LITERATURE REVIEW

A. Emotion Recognition in General

AI emotion recognition identifies and categorises human emotions from voice, writing, body language, and facial expressions. This technology lets robots comprehend and respond to human emotions, improving human-AI relations. Emotional intelligence improves decision-making and involvement in healthcare, education, customer service, and entertainment. Thus, emotion identification is crucial in these disciplines.

Traditional facial expression analysis has identified emotions like anger, grief, and happiness by detecting little mouth, eye, and eyebrow movements. Body language—posture, movement, and gestures—can also indicate emotion [3]. Modern AI systems, especially biometric ones, evaluate visual information to determine emotions. Recent research has focused on text and speech emotion extraction. This data is needed to interpret real-time human conversation in chatbots, virtual assistants, and call centres.

Natural Language Processing-based text-based emotion detection lets machines understand and categorise emotions from written content. Speech-based emotion recognition employs pitch and tone to evaluate the speaker's mood.

B. Emotion Recognition from Text

Emotion recognition from literature uses sentence or paragraph-level linguistic cues to extract emotions. Text-based emotion identification relies on word choice and sentence structure to convey emotions like joy, anger, fear, and sadness. Many models and methods are used for this.

➤ Techniques:

Natural Language Processing (NLP) is the foundation of emotion recognition from text. NLP techniques allow machines to interpret human language by analyzing syntax and semantics. Emotion detection from text relies on identifying words, phrases, or patterns associated with specific emotions. **Sentiment Analysis** focuses on determining the emotional polarity of a text, classifying it as positive, negative, or neutral [4]. Sentiment analysis, though often focused on the general sentiment of a text, can be expanded to identify specific emotions like fear or joy. **Lexical Analysis** involves identifying keywords or lexicons that indicate certain emotions. Predefined emotion lexicons, such as the NRC Emotion Lexicon, are commonly used to map words to specific emotional categories.

➤ Models

Bag of Words (BoW) model reduces text to words without considering order or grammar. Despite its simplicity, it may miss context, which is crucial for emotion identification.

- **TF-IDF:** Term Frequency-Inverse Document Frequency. The TF-IDF model statistically analyses word value in a corpus or document. It weights words by frequency to locate emotionally charged terms in a text.

Deep learning, especially LSTM and RNNs, has improved text emotion identification. Due to their sequential data handling, these models are good for capturing phrase contextual interactions.

Generational Pretrained Transformer (GPT) and Bidirectional Encoder Representations from Transformers (BERT) models are good at emotion recognition because they understand text context and semantics. Pre-trained on enormous datasets, these algorithms can accurately predict emotions by studying complex language patterns.

➤ Challenges

- Emotion recognition models often struggle to detect sarcasm, as the emotional tone implied by the words might differ from the literal meaning.
- Understanding the context of a conversation is essential for accurately identifying emotions. The same word or phrase can express different emotions depending on the context.
- Words or sentences can have multiple meanings or emotional connotations, making it difficult for models to assign the correct emotion.

C. Emotion Recognition from Speech

Emotion recognition from speech involves analyzing vocal cues such as tone, pitch, and rhythm to determine the speaker's emotional state. While text-based emotion detection focuses on linguistic content, speech-based systems primarily focus on **acoustic features** that reflect emotions [5]. The emotional state of a speaker can significantly affect the way they modulate their voice, making speech a rich source for emotion detection.

➤ Acoustic Features

Pitch may express many emotions. We usually mean eagerness or fury when we raise our voice and sadness or peace when we drop it. Speech tone reveals a person's emotions. Harsh or sharp tones can transmit negative feelings, whereas softer tones can convey positive ones. Impatience and anxiety may be reflected in speech rhythm and pace since people talk faster when worried and slower when delighted [6].

➤ Models

- **Hidden Markov Models (HMM)** are commonly used in speech emotion recognition due to their ability to model temporal sequences of acoustic features. They capture the dynamic aspects of speech, such as variations in tone or pitch, over time.
- **Support Vector Machines (SVM)** are effective in classifying emotions based on acoustic features by separating different emotional states into distinct categories.
- Deep learning models, particularly Convolutional Neural Networks (CNNs) and LSTMs, have shown great success in speech emotion recognition by learning complex

patterns in acoustic data and improving the accuracy of predictions [8].

➤ Challenges

- Variations in accents and dialects can affect speech patterns, making it challenging to accurately detect emotions across different speakers.
- Speech data is often collected in environments where background noise can interfere with the clarity of vocal features, impacting the performance of emotion recognition systems.
- The same acoustic features can be associated with multiple emotions, making it difficult to distinguish between closely related emotions like frustration and anger.

[9] created a multimodal emotion identification model to improve emotion detection. This model includes audio, video, and text. Their model shows that text and speech cues improve emotion categorisation, especially for complex emotions.

[10] showed that improved deep learning models for voice emotion identification can survive noise.

[11] introduced DeepMoji-like pre-trained models to text-based emotion recognition. Transfer learning from emoticon prediction improves DeepMoji's written content emotion recognition. By detecting subtle emotional cues in large social media datasets, our model enhanced emotion detection accuracy.

Deep learning and advanced acoustic feature extraction have enhanced spoken emotion recognition, whereas BERT and GPT have expanded text-based emotion recognition. These advances are enabling more complicated systems to precisely and instantly understand human emotions.

III. METHODOLOGY

This report uses secondary research to review past works on written and spoken emotion recognition. Academic papers, journal articles, and relevant research publications were reviewed to examine the field's diverse methods, concepts, and issues. The sources were chosen for their relevance, recentness, and capacity to illuminate audio analysis and NLP emotion recognition. IEEE, Google Scholar, and PubMed were used to research SVM, neural networks, and transformers.

IV. TECHNIQUES AND APPROACHES

A. Natural Language Processing (NLP) in Emotion Detection from Text

NLP is crucial to text emotion recognition. Natural Language Processing (NLP) lets computers understand, analyse, and respond to human speech. For emotion recognition, natural language processing (NLP) approaches analyse text data's sentiment, semantics, grammar, and other linguistic properties. The purpose is to detect emotional

words, phrases, and sentence patterns to extract their emotional essence through language intricacies.

Lexical analysis, sentiment analysis, and semantic modelling are NLP methods for text emotion identification. Sentiment analysis classifies text as good, negative, or neutral to determine its emotional tone [12].

Using lexicons or emotion dictionaries like the NRC Emotion Lexicon to categorise emotions like happiness, wrath, fear, and sadness might expand on this base level of emotion detection. Lexical analysis also deconstructs text into emotive words and phrases using rule-based systems or machine learning.

Semantic modelling improves natural language processing by emphasising word and phrase context. Transformers (e.g., BERT and GPT) have revolutionised natural language processing (NLP) by better understanding word relationships, enabling emotion recognition in complex or ambiguous texts.

NLP emotion recognition struggles with context-dependent expressions like sarcasm and ambiguity. Models must understand the text's tone and meaning as well as its words. Recent advances in contextualised word embeddings, deep learning, and BERT have helped machines better grasp emotions [13]. These models can detect subtle emotional overtones since they have been trained on enormous amounts of text and understand each word.

B. Acoustic Feature Extraction in Speech Emotion Detection

Acoustic properties like volume, pitch, tone, rhythm, and speed are used to distinguish spoken emotions. These voice features reveal an individual's emotions. A slower speech pace and softer tone may indicate melancholy or tranquilly, whereas a higher pitch and volume may indicate fury or agitation. These aural variations help identify emotions since they often convey more emotion than words.

Pitch is crucial to speech emotion recognition. Sound frequency is a good indicator of emotional arousal [14]. Joy, excitement, and nervousness are associated with louder voices, whereas sadness and disappointment with lower pitches. Tone and timbre indicate the speaker's mood and voice quality. Loudness and volume might indicate zeal, annoyance, or hatred.

Since emotions can alter speech, rhythm and tempo are also important. When concerned or irritated, people speak faster than when sad or pondering deeply.

These auditory data are extracted and analysed using specialised algorithms and machine learning. MFCCs, ZCR, and chroma feature extraction are used to analyse audio signals. These properties teach machine learning models to identify audio data patterns and assign emotions. Acoustic analysis can improve customer connections by detecting vocal contentment or discontent.

C. Machine Learning and Deep Learning Models for Emotion Recognition

Emotion identification systems that process spoken and written language use machine learning and deep learning. Data-driven emotion classification is possible with these models. Several models exist, each with pros and cons.

SVM is a common emotion recognition machine learning model because it categorises high-dimensional data well. SVMs can recognise emotion classes and create hyperplanes to partition them using acoustic data, making them perfect for speech emotion recognition. Computational constraints make it work well with simple datasets but poorly with large, complex ones.

Random Forest, another machine learning technique, is used extensively in emotion recognition [15]. It generates many decision trees and aggregates their results to improve classification accuracy. Random Forest can handle large datasets and complex feature interactions, making it a versatile text or speech emotion identification algorithm.

Deep learning models may handle unstructured data better than it, such as very semantically complicated language.

Convolutional Neural Networks (CNNs) are typically used to interpret images and videos, but their ability to handle structured data like audio signal spectrograms makes them effective at speech emotion recognition. CNNs can recognise local data patterns to identify emotion-related audio aspects. CNNs are useful for text emotion identification; they record local word associations by analysing text as word embeddings. A.I. and RNNs with LSTMs can recognise emotions in written and spoken language because they digest sequential material well [16].

These models understand context and time dynamics because they remember previous input (sentence phrases, speech fragments, etc.). LSTM models excel in identifying long-term dependencies in text or speech, and considering the whole sequence of inputs improves emotion recognition. Transformers' self-attention mechanisms have revolutionised text-based emotion identification.

These techniques let the model focus on key text elements. BERT and GPT, short for "Bidirectional Encoder Representations from Transformers" and "Generative Pre-trained Transformer," aim to identify emotions. BERT excels at emotion identification from complex or ambiguous texts because to its bidirectional context understanding [17].

Transformers are being studied for their greater efficiency and accuracy with sequential input, despite CNNs and RNNs' long-time supremacy in speech emotion identification.

Text and speech transformers are becoming more attractive to investigate due to their adaptability to unstructured data.

V. CHALLENGES IN EMOTION RECOGNITION

A. Multimodal Fusion

One of the major challenges in emotion recognition is the effective integration of multiple modalities, such as text and speech, to create a more comprehensive and accurate emotion recognition system. Multimodal fusion refers to the process of combining different types of data (e.g., linguistic and acoustic features) to enhance the accuracy and reliability of emotion detection. While each modality carries valuable emotional cues, merging them introduces several complexities.

Text data provides context and explicit emotional indicators through word choice and sentence structure, but it lacks the non-verbal cues available in speech, such as tone, pitch, and rhythm. On the other hand, speech data is rich in acoustic features that reveal the speaker's emotional state but often lacks sufficient semantic context to interpret the exact emotion being conveyed. The challenge lies in finding the right techniques to fuse these complementary modalities effectively [18]. For example, a sentence like "I'm fine" might indicate positive emotion in text, but when spoken with a sarcastic tone, it might convey the opposite emotion. Thus, multimodal fusion must take into account the temporal alignment of text and speech and determine which features to prioritize in each context.

Machine learning models used for multimodal fusion, such as early fusion, late fusion, or hybrid approaches, need to ensure that the unique features of each modality are adequately represented and interpreted. However, achieving this is complicated by differences in data types and structures. Early fusion involves combining raw features from both modalities before processing, but this can lead to information overload and inefficient training. Late fusion, on the other hand, processes each modality independently before merging their results, but it risks losing important correlations between modalities. Hybrid approaches attempt to balance these by aligning text and speech data at multiple stages in the model. Despite advancements, creating a seamless multimodal emotion recognition system remains a significant challenge due to the complexities of feature extraction, synchronization, and interpretation.

B. Ambiguity in Language and Speech

Language and voice ambiguity complicates emotion recognition. Cultural differences, personal experiences, and expression styles can greatly affect emotion perception and communication. Emotion recognition systems struggle to generalise across populations and circumstances due to this diversity [19].

Context, speaker intention, and culture can change the emotional connotations of a word or phrase in literature. For instance, "great" might mean happiness or sarcasm. Because they require deeper, frequently implied meanings, irony, sarcasm, and humour are hard for machines to identify. A phrase like "Oh, wonderful!" might represent excitement or irritation, which an emotion identification algorithm may miss without context.

Acoustic elements transmit emotions in speech, but voice pitch, tone, and speaking style can make them unclear. Some people speak with a monotonous voice, which may be misconstrued as a lack of emotion, while others show enthusiasm or fury with identical vocal signals, causing confusion [20]. Cultural influences also affect emotion expression and perception. Some cultures hide or quietly convey emotions like anger or sadness, while others express them openly. Cultural differences make it hard to create universal emotion recognition systems that work across languages and cultures.

C. Data Quality

An key challenge with emotion recognition is data quality. Emotion detection systems need enormous tagged data to train their machine learning models. Unlabelled writing or loud speech makes emotional data collection harder. Text-based emotion recognition is difficult when managing unstructured and unlabelled data like social media postings, emails, and conversations. Labelling such data with emotions is expensive and time-consuming. Human biases in

manually tagged data may also affect emotion recognition programs.

Noise, poor audio quality, and recording settings can greatly reduce speech emotion detection data quality. In loud environments like public spaces and phone calls, distortions and overlapping sounds make it harder to extract emotional indications like pitch and tone from voice data. Due to variances in acoustic properties induced by irregular recording equipment and conditions, training models that generalise across contexts is difficult.

Data quality also suffers from emotion parity issues. Joy and sadness are more common in training data than disgust and surprise. This mismatch can lead to biased algorithms that detect common emotions well but not rare ones. These data quality issues must be addressed to build reliable emotion identification algorithms that operate well in real life. To address data noise, labelling difficulties, and emotional imbalances, emotion detection systems need data augmentation, transfer learning, and unsupervised learning to improve accuracy and generalisability.

Table 1: Applications of Emotion Recognition

Application	Description
Customer Service & Chatbots	Enhances interactions by detecting emotional tone, allowing chatbots to respond empathetically and escalate issues if needed.
Healthcare	Monitors emotional states in patients to detect mental health issues like depression or anxiety, particularly in telemedicine.
Education	Tracks student engagement and emotional well-being to provide personalized learning and support teachers in addressing stress or frustration.
Entertainment & Gaming	Adapts game difficulty based on player emotions and enhances VR/AR experiences. Streaming services offer content recommendations based on emotional reactions.

VI. ETHICAL CONSIDERATIONS

Many ethical problems highlighted by emotion identification technology centre on privacy. Audio and text data collected to determine emotions reveal vulnerable mental health information without consent. Users may not understand their emotions are being tracked or analysed, therefore informed consent and data protection are crucial. Model bias can occur when emotion recognition algorithms are trained with gender, age, or culture-imbalanced datasets, resulting in biased or inaccurate results. A model built on Western data may misinterpret emotions in non-Western cultures, leading to misdiagnosis or misunderstandings. Similar to gender and age prejudices, different populations may receive different amounts of recognition. Misuse of this technology, especially for surveillance and emotion manipulation, is another problem. Concerns about governments or groups utilising emotion recognition to spy on or influence people raise questions of personal freedom, free expression, and the exploitation of such technologies in political campaigns or ads. For ethical reasons, emotion recognition technology must be open, have good data protection, and undergo continuous model development to be fair and accountable.

VII. FUTURE RESEARCH

Emotion identification will advance due to deep learning algorithms, which boost accuracy and adaptability. Modern transformer-based models (BERT, GPT) and other architectures are improving at identifying complex voice and text emotional cues. As multimodal systems become more common, emotion recognition becomes more complete. These systems use text, voice, body language, face expressions, and physiological sensors. These systems attempt to overcome the limitations of using one modality to understand human emotions by mixing input from many sources. Real-time emotion recognition is challenging, especially in loud or confusing conditions. Future systems could use edge computing and more efficient deep learning techniques for fast, accurate, and scalable emotion detection. These solutions may be useful in healthcare, interactive technology, and customer service. These innovations may revolutionise human-machine collaboration.

VIII. CONCLUSION

This article concluded that emotion recognition from text and speech is crucial to human-computer interaction. We addressed NLP, acoustic feature extraction, machine learning, and deep learning models that increase emotion recognition. These technologies have potential, but multimodal

integration, data uncertainty, and model bias remain. In AI applications like customer service, healthcare, and education, understanding human emotions can improve interactions and personalise them. Future advances in deep learning, multimodal systems, and real-time processing offer exciting prospects to overcome restrictions. Future research should reduce bias in models, improve multimodal integration for more accurate emotion recognition, and address ethical problems including privacy and exploitation of these technologies. Emotion recognition will evolve, helping AI systems become more compassionate and cleverer.

REFERENCES

- [1]. N. Alswaidan and M. E. B. Menai, "A survey of state-of-the-art approaches for emotion recognition in text," *Knowledge and Information Systems*, vol. 62, no. 8, pp. 2937-2987, 2020.
- [2]. N. Braunschweiler, R. Doddipatla, S. Keizer, and S. Stoyanchev, "Factors in emotion recognition with deep learning models using speech and text on multiple corpora," *IEEE Signal Processing Letters*, vol. 29, pp. 722-726, 2022.
- [3]. S. W. Byun, J. H. Kim, and S. P. Lee, "Multi-modal emotion recognition using speech features and text-embedding," *Applied Sciences*, vol. 11, no. 17, p. 7967, 2021.
- [4]. M. R. Makiuchi, K. Uto, and K. Shinoda, "Multimodal emotion recognition with high-level speech and text features," in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2021, pp. 350-357.
- [5]. K. Sailunaz, M. Dhaliwal, J. Rokne, and R. Alhajj, "Emotion detection from text and speech: A survey," *Social Network Analysis and Mining*, vol. 8, no. 1, p. 28, 2018.
- [6]. P. Singh, R. Srivastava, K. P. S. Rana, and V. Kumar, "A multimodal hierarchical approach to speech emotion recognition from audio and text," *Knowledge-Based Systems*, vol. 229, p. 107316, 2021.
- [7]. O. Verkholiyak, A. Dvoynikova, and A. Karpov, "A bimodal approach for speech emotion recognition using audio and text," *J. Internet Serv. Inf. Secur.*, vol. 11, no. 1, pp. 80-96, 2021.
- [8]. C. Wu, C. Huang, and H. Chen, "Text-independent speech emotion recognition using frequency adaptive features," *Multimedia Tools and Applications*, vol. 77, pp. 24353-24363, 2018.
- [9]. H. Xu, H. Zhang, K. Han, Y. Wang, Y. Peng, and X. Li, "Learning alignment for multimodal emotion recognition from speech," *arXiv preprint arXiv:1909.05645*, 2019.
- [10]. S. Yoon, S. Byun, and K. Jung, "Multimodal speech emotion recognition using audio and text," in *2018 IEEE Spoken Language Technology Workshop (SLT)*, 2018, pp. 112-118.
- [11]. S. K. Bharti, S. Varadhaganapathy, R. K. Gupta, P. K. Shukla, M. Bouye, S. K. Hingaa, and A. Mahmoud, "Text-based emotion recognition using deep learning approach," *Computational Intelligence and Neuroscience*, vol. 2022, p. 2645381, 2022.
- [12]. H. Lian, C. Lu, S. Li, Y. Zhao, C. Tang, and Y. Zong, "A survey of deep learning-based multimodal emotion recognition: Speech, text, and face," *Entropy*, vol. 25, no. 10, p. 1440, 2023.
- [13]. X. Cai, J. Yuan, R. Zheng, L. Huang, and K. Church, "Speech emotion recognition with multi-task learning," in *Interspeech*, 2021, pp. 4508-4512.
- [14]. B. T. Atmaja, K. Shirai, and M. Akagi, "Speech emotion recognition using speech feature and word embedding," in *2019 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2019, pp. 519-523.
- [15]. L. Pepino, P. Riera, L. Ferrer, and A. Gravano, "Fusion approaches for emotion recognition from speech using acoustic and text-based features," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6484-6488.
- [16]. M. S. Fahad, A. Ranjan, J. Yadav, and A. Deepak, "A survey of speech emotion recognition in natural environment," *Digital Signal Processing*, vol. 110, p. 102951, 2021.
- [17]. P. Hajek and M. Munk, "Speech emotion recognition and text sentiment analysis for financial distress prediction," *Neural Computing and Applications*, vol. 35, no. 29, pp. 21463-21477, 2023.
- [18]. T. Mittal, U. Bhattacharya, R. Chandra, A. Bera, and D. Manocha, "M3er: Multiplicative multimodal emotion recognition using facial, textual, and speech cues," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, vol. 34, no. 2, pp. 1359-1367.
- [19]. P. Kumar, V. Kaushik, and B. Raman, "Towards the explainability of multimodal speech emotion recognition," in *Interspeech*, 2021, pp. 1748-1752.
- [20]. X. Zhang, M. J. Wang, and X. D. Guo, "Multi-modal emotion recognition based on deep learning in speech, video and text," in *2020 IEEE 5th International Conference on Signal and Image Processing (ICSIP)*, 2020, pp. 328-333.