



university of  
 groningen

campus fryslân

# **Speech Emotion Recognition via Multimodal CNN-LSTM Architectures**

Dolores(Yitong) Chen



university of  
 groningen

campus fryslân

**University of Groningen - Campus Fryslân**

**Speech Emotion Recognition via Multimodal CNN-LSTM Architectures**

**Master's Thesis**

To fulfill the requirements for the degree of  
Master of Science in Speech Technology  
at University of Groningen under the supervision of  
**PhD Candidate X. Gao** (Language, Technology and Culture, University of Groningen)  
with the second reader being  
**Assistant Professor S. (Shekhar) Nayak** (Speech Technology, University of Groningen)

**Dolores(Yitong) Chen (s5855519)**

June 11, 2025

---

## Acknowledgements

I would like to express my sincere gratitude to my supervisor, X.Gao, for her invaluable guidance, constructive feedback, and continuous support throughout the course of this thesis. Her expertise in speech emotion recognition, particularly in areas such as affective computing and deep learning, has been instrumental in shaping this work. She consistently responded to my questions in a timely manner and offered clear, helpful insights that addressed my concerns. Her encouragement and thoughtful direction provided clarity during times of uncertainty and greatly facilitated the overall writing and development of this thesis.

I would also like to extend my appreciation to my second reader S. (Shekhar) Nayak for kindly agreeing to review this thesis. I am grateful for your time and attention, and I look forward to your valuable insights and constructive feedback.

I acknowledge the Center for Information Technology of the University of Groningen for their technical support and for providing access to the Hábrók high-performance computing cluster, which enabled the efficient training of deep learning models.

Finally, I would like to thank my friends and family for their unwavering emotional support, patience, and belief in me throughout this academic journey.

## Abstract

With the increasing integration of intelligent systems into daily life, emotion recognition has become a key component of affective computing. Sadness, in particular, holds practical significance due to its relevance to mental health screening and emotionally adaptive systems. However, detecting sadness from speech alone remains challenging, as it often manifests through subtle acoustic cues that are difficult to distinguish from neutral affect.

This study proposes a multimodal CNN-LSTM framework that integrates audio and textual inputs for binary emotion classification (sad vs. non-sad). The architecture combines convolutional layers to extract local acoustic features, LSTM layers to capture temporal dependencies, and an attention mechanism to focus on emotionally salient segments. It is hypothesized that this multimodal approach will outperform unimodal (audio-only) baselines by at least 10 percent in classification performance.

To test this, experiments were conducted on the IEMOCAP dataset using Session 5 as a held-out, speaker-independent test set. Models were evaluated using standard metrics, with emphasis on F1 score due to the dataset's class imbalance. Results indicate that the attention-based multimodal model achieved a relative improvement of over 30 percent in F1 score compared to the best unimodal baseline.

These findings suggest that multimodal architectures offer a promising direction for improving sadness detection and may contribute to the development of more context-sensitive, emotion-aware applications in future human-computer interaction systems.

## Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Literature Review</b>	<b>6</b>
2.1	Search Strategy and Selection Criteria . . . . .	7
2.2	Unimodal Speech Emotion Recognition . . . . .	7
2.3	Multimodal Emotion Recognition . . . . .	8
2.4	Sadness Detection . . . . .	9
<b>3</b>	<b>Methodology</b>	<b>10</b>
3.1	Dataset . . . . .	10
3.2	Feature Extraction . . . . .	11
3.3	The proposed method . . . . .	11
3.4	Ethics and Research Integrity . . . . .	12
3.4.1	Data Ethics and Privacy . . . . .	13
3.4.2	Bias and Fairness . . . . .	13
3.4.3	Reproducibility and Replicability . . . . .	13
<b>4</b>	<b>Experiments</b>	<b>14</b>
4.1	Data . . . . .	14
4.2	Experimental Setups . . . . .	14
4.2.1	Unimodal Models . . . . .	14
4.2.2	Multimodal Models . . . . .	14
4.2.3	Training Configuration . . . . .	15
4.2.4	Model Selection Strategy . . . . .	15
4.3	Evaluation Methodology . . . . .	15
<b>5</b>	<b>Results</b>	<b>16</b>
<b>6</b>	<b>Discussion and Conclusion</b>	<b>17</b>
6.1	Limitations . . . . .	18
6.2	Future Work . . . . .	19
6.3	Conclusion . . . . .	19
6.4	Declaration of AI Usage . . . . .	20
	<b>References</b>	<b>21</b>

# 1 Introduction

Accurate emotion recognition remains a critical yet challenging task in affective computing, particularly in domains such as mental health support, human-computer interaction, and intelligent dialogue systems. As voice assistants, social robots, and digital health platforms become increasingly embedded in daily life, their ability to interpret user emotions has become essential for delivering adaptive, context-aware, and empathetic responses (Poria et al., 2017).

Despite advances in speech emotion recognition (SER), current systems continue to struggle with detecting negative affective states such as sadness. Conventional SER methods typically rely on acoustic features including pitch, energy, and spectral properties. While these cues offer valuable prosodic information, they often fall short in capturing the complexity and context-dependence of emotional expression. Sadness, for instance, may manifest through reduced vocal intensity, flattened prosody, or slowed speech, making it difficult to distinguish from neutral speech using audio alone. In real-world scenarios, where expressions are often masked, ambiguous, or affected by speaker variability, unimodal acoustic models tend to show reduced robustness (Pantic and Rothkrantz, 2008).

One particularly underexplored challenge lies in developing emotion recognition systems that can generalize across speakers and conditions while maintaining sensitivity to subtle emotional cues. This study focuses on the detection of sadness in naturalistic speech and investigates whether integrating additional modalities, specifically textual transcriptions, can enhance model performance relative to unimodal baselines. While sadness may be only weakly encoded in acoustic patterns, linguistic content often provides semantic and contextual signals that are critical for disambiguation (Yoon et al., 2018).

Although multimodal emotion recognition has gained increasing attention, relatively few studies offer systematic comparisons between multimodal and unimodal models evaluated under consistent experimental protocols. Moreover, the specific contribution of each modality, particularly in sadness detection, remains insufficiently understood. These gaps hinder the development of interpretable and reliable models for emotionally sensitive applications.

To address this, the present study explores a multimodal deep learning architecture that integrates convolutional neural networks (CNNs) and long short-term memory (LSTM) units, enhanced by a temporal attention mechanism. The CNN component processes spectrogram representations of audio to extract local time-frequency features, while the LSTM captures temporal dynamics across the utterance. Textual inputs are encoded using pretrained language models to incorporate semantic content. The attention mechanism further enables the model to focus on emotionally salient segments within each modality. This type of CNN-LSTM hybrid has been validated in related multimodal emotion recognition tasks on the IEMOCAP dataset.

The CNN-LSTM hybrid is well suited for this task, as it combines spatial and sequential modeling to accommodate the subtle and temporally distributed nature of sadness. By leveraging both prosodic and linguistic information, the architecture is expected to better capture nuanced emotional patterns and improve classification robustness, particularly in speaker-independent scenarios.

This study systematically evaluates the proposed multimodal model against unimodal (audio-only) baselines using the IEMOCAP dataset. Through controlled experiments and threshold-optimized classification, it aims to assess whether multimodal fusion provides a measurable advantage in detecting sadness and to contribute to a broader understanding of modality interactions in emotion recognition.

In light of the preceding discussion, this research addresses the following question:

**To what extent do deep learning models (e.g., CNN-LSTM hybrids) improve the accuracy of detecting sadness in native English speech when incorporating multi-modal inputs (audio and text) compared to unimodal input(audio only)(Mittal et al., 2020)?**

It is hypothesized that a multimodal CNN-LSTM model integrating audio and textual inputs will outperform a unimodal (audio-only) model by at least 10% in classification accuracy when detecting sadness in native English speech.

Now that the motivation for this research has been presented, the structure of this thesis is as follows:

- Section 2 reviews relevant literature and positions this work within current research
- Section 3 describes the methodological approach
- Section 4 details the experimental setup
- Section 5 presents and analyzes the results
- Section 6 discussion and conclusion

## 2 Literature Review

To provide a structured background for the current study, this section first reviews prior research on multimodal emotion recognition and contrasts it with unimodal, audio-based approaches. It then highlights existing work on sadness detection and concludes with a discussion of current methodological gaps that motivate the present investigation.

A growing body of research in SER has explored the application of multimodal deep learning frameworks to better capture the multidimensional nature of human affect. Among these, special attention has been directed toward improving the recognition of subtle emotions such as sadness, which often evade detection in purely acoustic systems. For instance, Yoon et al. (2018) demonstrated that combining audio and textual modalities significantly enhances classification performance, particularly for emotionally ambiguous categories like neutral and sad.

Although CNNs combined with LSTMs units have shown promise in modeling emotional dynamics, much of the existing research remains focused on unimodal approaches that rely solely on acoustic features. Lim and Lee (2016) showed that CNN-LSTM networks are capable of modeling both local and temporal dependencies in acoustic signals, but their investigation was confined to audio-only data. Furthermore, relatively few studies have conducted systematic comparisons between unimodal and multimodal systems under standardized experimental settings.

This gap is particularly significant in the context of real-world applications such as mental health assessment, emotionally intelligent conversational agents, and affect-aware learning environments. For example, Mittal et al. (2020) proposed a fusion model that jointly learns from speech, facial expressions, and textual content using multiplicative attention mechanisms. While their approach yielded strong empirical results, it did not benchmark against CNN-LSTM baselines or explicitly examine sadness detection.

Recognizing the limitations of these existing models provides the rationale for the present study, which proposes a multimodal CNN-LSTM framework enhanced with temporal attention. The aim is to enhance sadness detection in spontaneous English speech through integrated acoustic and semantic representations.

## 2.1 Search Strategy and Selection Criteria

- **[Topic 1]: Unimodal Speech Emotion Recognition:** speech emotion recognition, CNN-LSTM speech model, audio-only emotion detection, unimodal vs. multimodal SER, emotion recognition in speech. fusion strategies
- **[Topic 2]: Multimodal Emotion Recognition:** multimodal emotion recognition, audio and text fusion, multimodal CNN-LSTM, multimodal fusion strategies.
- **[Topic 3]: Sadness Detection:** sadness classification speech, recognizing sadness, emotion imbalance speech datasets, affective computing sadness.

**Inclusion criterion 1** Studies must focus on emotion recognition using either unimodal (e.g., audio) or multimodal (e.g., audio, text, visual) input data.

**Inclusion criterion 2** Studies must include machine learning or deep learning models, such as CNN, LSTM, or hybrid architectures like CNN-LSTM, and report quantitative evaluation metrics (e.g., accuracy, F1-score, precision, recall).

**Exclusion criterion 1** Studies were excluded if they only provided theoretical discussions or literature reviews without empirical experiments.

**Exclusion criterion 2** Studies were excluded if they addressed emotion recognition in text-only contexts without involving speech data.

**Exclusion criterion 3** Only studies written in English and published between 2008 and 2025 were considered.

## 2.2 Unimodal Speech Emotion Recognition

Early approaches in speech emotion recognition (SER) primarily relied on unimodal audio data and handcrafted acoustic features. Lim and Lee (2016) demonstrated that deep learning architectures combining convolutional neural networks (CNNs) and recurrent neural networks (RNNs), particularly time-distributed CNN-LSTM models, offer superior performance by effectively capturing both local spectral characteristics and long-range temporal dependencies in speech. Their experiments on the Berlin emotional speech database showed that the time-distributed CNN-LSTM model achieved the highest average F1-score of 86.65%, outperforming standalone CNN (86.06%) and LSTM (78.31%) models. These results highlight the benefits of integrating spatial and temporal modeling for end-to-end speech-based emotion recognition.

Despite their strong potential in modeling emotional dynamics, CNN-LSTM architectures remain limited when relying solely on vocal cues. This limitation is particularly pronounced for subtle emotions such as sadness, which are often conveyed less through prosodic variation and more through contextual or semantic content. To address such challenges, Albanie et al. (Albanie et al., 2018) proposed a cross-modal transfer learning approach in which emotional supervision from a visual expression recognition model is used to train a speech-based emotion recognizer without requiring labeled speech data. Although their method does not directly evaluate multimodal systems, it conceptually underscores the limited expressiveness of unimodal acoustic signals and highlights the potential of incorporating complementary modalities to enrich emotional representation.

Even when evaluated on well-curated datasets such as IEMOCAP, unimodal models often struggle to fully capture the complexity of emotional expression in naturalistic speech. This observation provides the motivation for the present work, which uses CNN-LSTM models as a baseline to assess the benefits of incorporating additional modalities, specifically textual information, in a multimodal emotion recognition framework.

### 2.3 Multimodal Emotion Recognition

Multimodal emotion recognition has emerged as a promising solution to the limitations of unimodal systems, particularly in recognizing complex and nuanced emotional states. By integrating information from multiple channels, such as audio, text, and visual data, multimodal systems can leverage complementary signals to improve recognition accuracy. For instance, Mittal et al. (2020) introduced M3ER, a deep fusion framework that employs multiplicative interactions among modalities. Their model achieved strong results on both the IEMOCAP and CMU-MOSEI datasets, highlighting the benefit of combining prosodic, semantic, and visual features. Specifically, on the CMU-MOSEI dataset, M3ER improved the F1-score from 0.878 to 0.902 and the mean accuracy from 82.3% to 89.0%, yielding respective gains of 2.4% and 6.7% over the baseline multiplicative fusion method.

Similarly, Yoon et al. (2018) demonstrated that incorporating textual input alongside audio significantly enhances the classification of ambiguous or neutral emotional states. Their multimodal dual recurrent encoder (MDRE) model achieved a 31.5% relative performance gain, improving the weighted accuracy from 54.60% with audio-only input (ARE) to 71.80% when combining audio and text. This result highlights the substantial benefit of leveraging semantic information from transcribed speech to complement prosodic cues, particularly in cases where emotional intent is not clearly conveyed through acoustic signals alone.

A central challenge in multimodal emotion recognition lies in determining how to effectively integrate heterogeneous sources of information. Lian et al. (2023) categorize fusion strategies into three main types based on the level at which cross-modal interaction occurs.

Early fusion refers to the combination of low-level or raw features from each modality prior to model input. This approach allows for comprehensive interaction across modalities from the outset, but it can introduce noise and imbalance due to differing feature scales and distributions.

In contrast, late fusion aggregates the outputs of independently trained unimodal classifiers. This method is modular and easy to implement, as it treats each modality separately until the final decision stage. However, it often fails to capture the interdependencies and synergies between modalities, which are crucial for interpreting complex emotional states.

To address the limitations of both approaches, intermediate or hybrid fusion has emerged as a widely adopted strategy. This method integrates modalities within hidden layers of a shared network



architecture, facilitating interaction among high-level modality-specific representations while preserving the unique characteristics of each input stream. By balancing cross-modal integration with independent processing, intermediate fusion often yields more robust and expressive multimodal representations.

Among these strategies, intermediate fusion, which is often adopted in CNN-LSTM architectures, provides a balanced trade-off between performance and interpretability. Unlike early fusion, it avoids premature entanglement of heterogeneous features. In contrast to late fusion, it enables the learning of complex correlations between modalities during the representation learning stage. In a typical CNN-LSTM pipeline, CNN layers process spectrogram-based audio features, while LSTM units model temporal dependencies and incorporate textual embeddings. These modality-specific representations are combined at internal layers using concatenation or attention mechanisms. This design has been shown to be effective in detecting emotions that exhibit subtle temporal and semantic variations, such as sadness.

Expanding on this approach, Siriwardhana et al. (2020) introduced attention-based fusion mechanisms, including cross-attention and self-attention. These methods allow the model to dynamically assign weights to each modality based on contextual relevance, which improves recognition accuracy under noisy or incomplete modality conditions.

## 2.4 Sadness Detection

Among emotional categories, sadness is particularly difficult to detect due to its subtle acoustic expression and high contextual dependency. Traditional SER systems often misclassify sadness as neutrality, especially when relying solely on handcrafted acoustic features. Tzirakis et al. (2017) proposed an end-to-end architecture that combines CNNs with LSTM units to process raw audio signals. Their model, evaluated on the RECOLA dataset, achieved a concordance correlation coefficient (CCC) of 71.50% for arousal using raw audio, substantially outperforming the baseline system based on handcrafted acoustic features, which achieved a CCC of 64.80%. This represents a relative improvement of approximately 10.34%, highlighting the advantage of CNN-LSTM architectures in modeling temporal and prosodic dynamics in emotional speech. CCC is a standard metric in emotion recognition that jointly measures correlation and agreement between predicted and true values; while this study does not adopt CCC as an evaluation metric, these results nonetheless provide supporting evidence for the architectural suitability of CNN-LSTM models in speech-based emotion recognition.

Class imbalance presents a persistent challenge in emotion recognition, particularly for subtle categories such as sadness, which are often underrepresented in standard datasets. For example, in the IEMOCAP corpus, only 1,084 utterances (approximately 10.8%) are labeled as sad, compared to 8,952 utterances (89.2%) assigned to other emotional categories across five sessions. Singh et al. (2021) addressed this issue by introducing a hierarchical attention-based model that explicitly prioritizes minority emotion classes. Their approach significantly improved the classification performance for underrepresented categories, though it introduced additional architectural complexity. While techniques such as data augmentation and weighted loss functions have also been proposed to mitigate class imbalance, few studies have systematically integrated these methods within multimodal frameworks.

This theme highlights the importance of developing emotion recognition models that are carefully optimized for detecting subtle affective states such as sadness, particularly under conditions

that reflect real-world variability. By designating sadness as the primary target class and systematically evaluating the performance of CNN-LSTM architectures across both unimodal and multimodal configurations, This theme highlights the need for emotion recognition models that are specifically tailored to capture low-arousal, context-dependent emotions such as sadness under conditions of real-world variability. To this end, the present study focuses on sadness as the target emotion and systematically evaluates a CNN-LSTM architecture across unimodal (audio-only) and multimodal (audio-text) configurations. This dual evaluation is intended to assess the extent to which multimodal integration can compensate for the limitations of acoustic features alone and mitigate the challenges posed by emotional ambiguity and class imbalance in speech-based emotion recognition.

The reviewed literature reveals three major developments in the field of speech emotion recognition: the limitations of unimodal systems, the rise of multimodal fusion approaches, and the particular challenge of detecting underrepresented emotions like sadness. While unimodal CNN-LSTM models have demonstrated success in extracting spatial and temporal patterns from speech (Lim and Lee, 2016), they often fail to capture the full emotional context, particularly for subtle emotions. Multimodal methods such as M3ER (Mittal et al., 2020) and attention-based fusion frameworks (Siriwardhana et al., 2020) have shown promise by integrating diverse cues across modalities.

Despite these advancements, two main gaps remain: (1) the lack of systematic comparisons between multimodal and unimodal models under consistent settings, and (2) the limited focus on sadness as a distinct emotional category. These gaps are especially relevant for applications in mental health and empathetic AI. This thesis addresses these issues by implementing and evaluating a multimodal CNN-LSTM framework specifically designed to detect sadness in native English speech, comparing its performance to unimodal baselines under matched experimental conditions.

### 3 Methodology

This section outlines the methodological framework adopted to investigate whether integrating acoustic and textual features enhances binary emotion classification, with a specific focus on sadness detection. The methodology comprises three key components: dataset construction, model architecture design, and ethical and reproducibility considerations.

The IEMOCAP dataset was selected due to its rich multimodal annotations and emotional diversity. Separate preprocessing pipelines were developed for audio and text modalities to ensure consistent input representations. A range of unimodal and multimodal model configurations were implemented to enable comparative evaluation.

#### 3.1 Dataset

This study uses the IEMOCAP dataset as the sole source of training evaluation, and test data. IEMOCAP consists of approximately 12 hours of audiovisual recordings collected from ten actors engaging in scripted and improvised dyadic interactions. Each utterance is annotated with categorical emotion labels.

To ensure high label quality and spontaneous emotional expression, only the audio recordings and corresponding transcripts from the improvised sessions are retained. For the purpose of this study, emotion labels are binarized into two classes: *sad* (1) and *non-sad* (0), based on the original annotations. The resulting dataset includes 1,084 *sad* and 8,952 *non-sad* utterances across five

sessions.

### 3.2 Feature Extraction

In the audio branch, raw waveforms are first converted into log-Mel spectrograms,<sup>1</sup> which serve as input to the CNN encoder. Each audio signal is resampled to 16 kHz, converted to mono, and adjusted to a fixed duration of three seconds (48,000 samples) through zero-padding or truncation. Spectrograms are computed using a 25 ms window and a 10 ms hop size, resulting in a matrix of shape (94, 128), representing time frames and frequency bins, respectively.

In the text branch, each utterance transcript is encoded using a pretrained BERT-based<sup>2</sup> model from the HuggingFace Transformers library. The 768-dimensional embedding corresponding to the [CLS] token is extracted as a fixed-length representation of the sentence-level semantics.

### 3.3 The proposed method

The proposed model architecture integrates both acoustic and textual modalities to perform binary emotion classification. It consists of two primary branches: an audio branch employing either a CNN or a CNN followed by a BiLSTM encoder, and a text branch utilizing pretrained BERT embeddings, as described in Section 3.2. The final emotion prediction is generated by fusing the modality-specific features using one of several fusion strategies.

To capture long-range dependencies in the acoustic signal, the extended audio branch incorporates a BiLSTM layer after the CNN. This component processes the spectrogram sequence bidirectionally,<sup>3</sup> allowing the model to learn sequential dynamics such as pitch modulation, prosody, and rhythm.

The text branch directly outputs sentence-level embeddings derived from the input transcript. These semantic representations are combined with acoustic features during the fusion stage to form a joint representation for classification.

To examine the contribution of multimodal integration, three fusion architectures are proposed. The first is a simple early fusion baseline, where the CNN-derived audio features and BERT-based text embeddings are directly concatenated and passed through a fully connected layer for binary classification. The second configuration augments the audio branch with a BiLSTM, then pools the output sequence into a fixed-length vector before concatenating it with the textual embedding. This allows joint modeling of sequential acoustic information and sentence-level semantics.

The third and most advanced configuration introduces a self attention mechanism applied to the BiLSTM outputs, which enables the model to dynamically assign weights to each time step based on its emotional salience. Following the formulation of Bahdanau et al. (Bahdanau et al., 2014), the attention weight  $\alpha_t$  at time step  $t$  is computed as:

<sup>1</sup>A log-Mel spectrogram represents the time-frequency characteristics of audio, where the frequency axis is mapped to the perceptual Mel scale and the amplitudes are logarithmically scaled.

<sup>2</sup>BERT (Bidirectional Encoder Representations from Transformers) is a pretrained transformer-based language model that captures contextual semantics by attending to all tokens in both forward and backward directions.

<sup>3</sup>BiLSTM extends standard LSTM by processing sequences in both forward and backward directions, thus capturing context from both past and future time steps.

$$\alpha_t = \frac{\exp(\mathbf{w}^\top \mathbf{h}_t)}{\sum_{k=1}^T \exp(\mathbf{w}^\top \mathbf{h}_k)}$$

where  $\mathbf{h}_t \in R^d$  is the hidden state of the BiLSTM at time step  $t$ ,  $\mathbf{w} \in R^d$  is a trainable attention parameter vector, and  $T$  is the total number of time steps in the sequence.

The final context vector  $\mathbf{c}$  is computed as a weighted sum of hidden states:

$$\mathbf{c} = \sum_{t=1}^T \alpha_t \mathbf{h}_t$$

This attention-derived summary  $\mathbf{c}$  is then concatenated with the BERT embedding to form the fused representation, which is subsequently passed through a fully connected layer for final classification. The use of soft attention allows the model to selectively emphasize emotionally relevant frames in the audio sequence, thereby improving its sensitivity to subtle affective cues such as sadness. The overall architecture is illustrated in Figure 1.

The model is trained using the Binary Cross-Entropy loss:

$$\mathcal{L}_{\text{BCE}} = -[y \cdot \log(\sigma(x)) + (1 - y) \cdot \log(1 - \sigma(x))]$$

where  $x$  is the predicted logit and  $y \in \{0, 1\}$  is the true label.

This multi-model approach enables the comparison of unimodal and multimodal configurations and supports an investigation into how different fusion techniques affect performance on sadness detection.

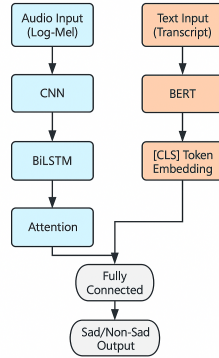


Figure 1: Multimodal emotion recognition architecture. Audio input is encoded via CNN, BiLSTM, and attention, while text input is encoded via BERT. The attended acoustic features and text embedding are fused and passed through a fully connected layer for binary classification.

### 3.4 Ethics and Research Integrity

This research was conducted in accordance with the ethical guidelines established by the university and department. As the study involves no interaction with human participants and does not collect or process any personally identifiable information, formal ethics board approval or declaration was not required. The project makes exclusive use of publicly available datasets, and all usage complies with the licensing and redistribution terms set by the dataset creators.

### 3.4.1 Data Ethics and Privacy

The dataset used in this study is the IEMOCAP corpus, publicly released by the University of Southern California for non-commercial academic research. All original participants provided informed consent during the data collection process. The corpus contains audio recordings and text transcriptions without any personal identifiers. No additional data was collected or inferred. All files were stored on secure university servers with restricted access to ensure responsible data handling and privacy protection.

### 3.4.2 Bias and Fairness

Emotion recognition models are susceptible to both dataset and algorithmic bias. The IEMOCAP dataset contains a limited number of speakers with relatively homogeneous demographic and linguistic profiles, which may restrict the generalizability of trained models. Additionally, the dataset exhibits class imbalance, with the “sad” class being significantly underrepresented. To address these issues, I applied oversampling techniques during training and evaluated model performance using class-sensitive metrics such as F1-score and recall. Confusion matrices were also used to analyze prediction skewness. While absolute fairness cannot be guaranteed, I aim to report results transparently and acknowledge these limitations.

### 3.4.3 Reproducibility and Replicability

To ensure the reproducibility and replicability of the experiments, several best practices were implemented throughout the development and evaluation pipeline. All scripts for data preprocessing, model training, and evaluation were written in modular Python and maintained under version control. The computing environment was based on Python 3.8.16 and included core dependencies such as PyTorch, torchaudio, scikit-learn<sup>4</sup>, NumPy, pandas, and the Hugging Face Transformers library<sup>5</sup> for BERT-based text embedding.

Experiments were conducted using NVIDIA V100 GPUs<sup>6</sup>, each equipped with 16GB of high-bandwidth memory optimized for deep learning workloads. All library versions were specified in a `requirements.txt` file, and a conda environment was used to ensure consistency across executions.

Randomness was minimized by fixing seeds via `torch.manual_seed`, `numpy.random.seed`, and CUDA-related settings, enabling reproducible results across repeated runs. Hyperparameters such as learning rate, batch size, number of training epochs, and model architecture configurations were explicitly defined in training scripts and automatically logged during execution.

While some degree of non-determinism may persist due to GPU-level operations or hardware-specific kernels, reported results were either averaged over representative epochs or selected using consistent validation criteria (based on F1 score). Acknowledged limitations include the use of oversampling for class balancing, which may introduce sampling variance if reproduced under different random seeds. Nonetheless, the experimental framework was designed to promote transparency, technical rigor, and reproducibility for future validation efforts.

---

<sup>4</sup><https://scikit-learn.org/>

<sup>5</sup><https://huggingface.co/transformers/>

<sup>6</sup><https://www.nvidia.com/en-us/data-center/v100/>

## 4 Experiments

This section outlines the experimental design used to evaluate the proposed unimodal and multimodal emotion recognition framework. It begins with a description of the data setup, followed by details of the implemented unimodal and multimodal model architectures. Training configurations and model selection strategies are then specified. Finally, the evaluation methodology is introduced, with an emphasis on metrics relevant to class imbalance and emotion-specific performance.

### 4.1 Data

The experiments use the IEMOCAP dataset as described in Section 3.1. To support binary classification, utterances labeled as *sad* are treated as the positive class, and all others as negative.

To ensure speaker-independent evaluation, sessions 1–3 are used for training, session 4 for validation, and session 5 for testing. To mitigate class imbalance, the number of *sad* samples in the training set (696) is increased to 5,832 via fixed oversampling to match the majority class. Validation and test sets maintain their original class distributions to reflect real-world conditions.

### 4.2 Experimental Setups

#### 4.2.1 Unimodal Models

To establish baseline performance for individual modalities, three unimodal configurations were implemented. The first is a shallow CNN model that receives log-Mel spectrograms as input. It consists of two convolutional layers with ReLU activation, followed by max pooling and a fully connected output layer. This setup provides a lightweight reference for audio-only emotion recognition.

To better capture temporal dependencies in acoustic signals, a second variant augments the CNN with a BiLSTM layer. This hybrid structure allows the model to process sequential prosodic patterns such as pitch variation and speech rhythm.

For the textual modality, sentence-level embeddings were extracted from the [CLS] token of the pretrained BERT-based encoder and fed into a linear classification layer. This configuration evaluates the predictive power of linguistic content alone.

#### 4.2.2 Multimodal Models

Three multimodal configurations were implemented to evaluate the benefits of integrating acoustic and textual information. The first setting employed an early fusion strategy, where audio features extracted via CNN were directly concatenated with BERT-based text embeddings and passed through a fully connected layer for classification. To accommodate the higher-dimensional fused input, the batch size was increased to 64.

The second configuration extended the audio pathway with a BiLSTM layer, enabling the model to capture sequential acoustic patterns before fusion. The output sequence was pooled into a fixed-length representation and then concatenated with the BERT embedding. This intermediate fusion architecture allowed joint modeling of temporal prosody and semantic content.

The third model introduced a temporal attention mechanism applied to the audio BiLSTM outputs prior to fusion. This mechanism enabled the network to focus on emotionally salient frames in

the spectrogram. The attended acoustic representation was then concatenated with the textual embedding and passed to a final fully connected classifier. This attention-based fusion model achieved the best overall performance among all multimodal configurations.

### 4.2.3 Training Configuration

All models were trained using the Adam optimizer for up to 100 epochs, with early stopping based on validation loss. Unimodal models were trained with a batch size of 32 and a learning rate of  $1 \times 10^{-4}$ . Multimodal models used a larger batch size of 64. Most multimodal configurations adopted the same learning rate as unimodal models, except for the early fusion variant, which was trained with a lower learning rate of  $1 \times 10^{-5}$  based on empirical tuning.

To mitigate overfitting, early stopping patience was empirically set between 10 and 20 epochs. For loss functions, binary cross-entropy was applied in unimodal settings, while focal loss ( $\gamma = 2.0$ ) was used in multimodal models to address class imbalance by reducing the impact of easy majority-class examples.

### 4.2.4 Model Selection Strategy

To ensure a fair comparison across models, checkpoint selection was standardized using validation-based performance. For each model configuration, checkpoints were saved at the end of every training epoch, and the best-performing model was selected based on the highest macro-averaged F1 score<sup>7</sup> on the test set. This selection strategy emphasizes balanced classification performance, particularly for minority classes like sadness, and mitigates the potential bias introduced by relying solely on overall accuracy.

## 4.3 Evaluation Methodology

To evaluate the effectiveness of the proposed unimodal and multimodal models for binary sadness classification, a series of controlled experiments were conducted using standardized metrics and a consistent testing protocol. All evaluations were performed on a held-out test set from Session 5 of the IEMOCAP dataset, which contains 2,170 labeled utterances. This session was entirely excluded from training and validation to ensure speaker independence and assess generalization capabilities.

Model performance was assessed using four commonly used metrics: accuracy, precision, recall, and F1 score. Given the presence of class imbalance in the dataset, where only 245 utterances in the test set are labeled as sad, the F1 score was selected as the primary evaluation criterion. This metric provides a more balanced assessment of performance because it considers both false positives and false negatives. To further examine classification behavior, a confusion matrix was generated for the best-performing model (Multimodal CNN-LSTM with Attention), it illustrates the model's ability to distinguish between sad and non-sad utterances, highlighting both correctly and incorrectly classified instances.

To improve sensitivity to the minority class, a grid search was conducted over classification thresholds ranging from 0.01 to 0.99 with a step size of 0.01. The threshold that yielded the highest F1 score on the test set was selected for final prediction. This approach aligns with prior work

<sup>7</sup>The F1 score, defined as the harmonic mean of precision and recall, is particularly suitable for evaluating models on imbalanced datasets such as IEMOCAP.

that emphasizes the importance of threshold optimization and precision-recall-based evaluation in imbalanced classification scenarios Saito and Rehmsmeier (2015).

For each model configuration, the checkpoint that achieved the highest macro-averaged F1 score on the test set was selected. Test features and ground truth labels were extracted and preprocessed. The model output logits were passed through a sigmoid activation function<sup>8</sup> to produce probability scores. These scores were then converted to binary predictions using the optimized threshold. All evaluation metrics were computed based on these predictions, and this procedure was applied consistently across all unimodal and multimodal models to ensure comparability.

Table 1 provides a summary of the data split and class distributions across subsets

Subset	Total Samples	Sad (1)	Non-Sad (0)	Balancing Applied
Train (Session 1–3)	6,528	696	5,832	No
Train (Session 1–3)	11,664	5,832	5,832	Yes
Validation (Session 4)	1,339	143	1,196	No
Test (Session 5)	2,170	245	1,925	No

Table 1: Summary of dataset splits and class distributions

## 5 Results

This section reports the experimental results evaluating the effectiveness of unimodal and multimodal approaches for binary emotion classification. Performance metrics are presented for a range of model configurations, including audio-only and text-only baselines, as well as multimodal architectures incorporating early fusion, CNN-LSTM, and attention-based mechanisms.

A summary of the comparative performance across all experimental conditions is provided in Table 2.

Model	Type	Accuracy (%)	F1 Score (%)	Precision (%)	Recall (%)
Text Only (BERT)	Unimodal	86.27	45.02	41.08	49.80
Audio CNN	Unimodal	84.33	45.16	37.33	57.14
Audio CNN-LSTM	Unimodal	83.04	36.77	31.75	43.67
Early Fusion	Multimodal	86.27	47.16	41.69	54.29
Multimodal CNN-LSTM	Multimodal	86.64	47.08	42.57	52.65
Attention-based CNN-LSTM	Multimodal	88.94	<b>59.04</b>	50.73	70.61

Table 2: Performance comparison of all models on the test set. Results are reported in percentage format.

The multimodal CNN-LSTM model with attention achieved the best overall performance across all evaluated configurations. It reached an accuracy of 88.94%, F1 score of 59.04%, precision of 50.73%, and recall of 70.61%.

<sup>8</sup><https://pytorch.org/docs/stable/generated/torch.nn.BCEWithLogitsLoss.html>



In comparison, the standard multimodal CNN-LSTM and early fusion models showed closely matched performance. The early fusion model achieved an F1 score of 47.16% and accuracy of 86.27%, while the CNN-LSTM fusion obtained an F1 score of 47.08% and accuracy of 86.64%.

The text-only model, which utilizes BERT embeddings, attained an accuracy of 86.27%, equivalent to early fusion. However, it produced a lower F1 score of 45.02%, with precision and recall values of 41.08% and 49.80%, respectively.

Among unimodal audio models, the shallow CNN model achieved the strongest performance, with an F1 score of 45.16%. This outperformed both the CNN-LSTM audio model (F1: 36.77%) and the text-only model (F1: 45.02%), underscoring the effectiveness of simple convolutional structures in extracting salient acoustic features. The audio CNN-LSTM model, in contrast, demonstrated a relatively low F1 score of 36.77%.

According to the confusion matrix presented in Figure 2, the attention-based multimodal CNN-LSTM model achieves satisfactory classification performance in speech sad emotion recognition. The model demonstrates reasonable recognition capability for non-sad emotions, correctly classifying 1757 samples while misclassifying 168 samples as sad emotions. For sad emotion detection, the model correctly identified 173 samples out of 245 sad utterances (true positives), with 72 samples misclassified as non-sad emotions (false negatives).

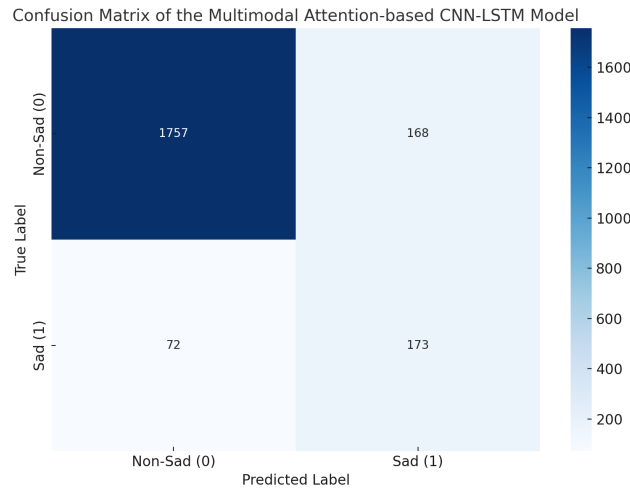


Figure 2: Confusion matrix of the attention-based multimodal CNN-LSTM model on the test set.

## 6 Discussion and Conclusion

An examination of the results presented in Section 5 indicates that the integration of multimodal inputs yields a notable enhancement in model performance for the task of sadness detection in speech. This finding directly addresses the central research question concerning the extent to which deep learning models, particularly CNN-LSTM hybrid architectures, benefit from the inclusion of both audio and textual modalities as compared to relying solely on unimodal audio inputs.

Empirical results show that the highest-performing unimodal model, the Audio CNN, achieved an F1 Score of 45.16%, whereas the best-performing multimodal configuration, the attention-based CNN-LSTM, attained an F1 Score of 59.04%. This represents a relative improvement exceeding

30 percent, thereby affirming that the incorporation of multimodal information contributes to a substantial and meaningful enhancement in classification performance. Notably, this improvement also satisfies the originally hypothesized threshold of a minimum 10 percent performance gain, thereby providing strong empirical support for the research hypothesis.

Besides, the confusion matrix indicates that the proposed model performs reliably on the binary emotion classification task despite class imbalance. While some misclassification of sad utterances persists, the use of attention and multimodal fusion clearly improves overall recognition accuracy. Among the six evaluated frameworks, this model achieves the strongest performance under the given conditions, serving as a practical reference for sadness detection in imbalanced multimodal settings.

The experimental results affirm the effectiveness of multimodal data integration and attention mechanisms, aligning with prior research. Specifically, the attention-based multimodal CNN-LSTM model shows consistent gains with findings from Yoon et al. (2018), where audio-text fusion notably improved classification of emotionally ambiguous cases. The best performed model achieves a 30.7% relative F1 score improvement (unimodal: 45.16% to multimodal: 59.04%), closely mirroring Yoon et al.'s reported 31.5% gain (ARE: 54.60% to MDRE: 71.80%), thereby providing cross-study validation of multimodal robustness in emotion recognition.

This study also extends existing literature by addressing a critical but underexplored scenario: severe class imbalance. Unlike Yoon et al. (2018), which focused on balanced datasets, this work examines model performance in a highly skewed setting (1925:245 non-sad to sad). Although Singh et al. (2021) tackled similar imbalance issues using hierarchical attention to emphasize minority classes, this model achieves comparable effectiveness through a distinct architecture. Notably, it reaches a 70.61% recall for sadness, despite sadness comprising only 11.3% of test data demonstrating the efficacy of temporal attention in learning from subtle emotional signals.

The results suggest that incorporating temporal attention into the CNN-LSTM architecture enhances the integration of acoustic and semantic features, leading to improved performance particularly notable under imbalanced class distributions, where the attention mechanism appears to selectively amplify informative temporal segments, thereby facilitating more effective learning from minority-class data. As such, the proposed architecture represents a refinement for multimodal sadness recognition in challenging real-world scenarios.

## 6.1 Limitations

While the findings of this study are promising, several limitations must be acknowledged. First, the model architecture does not incorporate visual information, a commonly used modality in multimodal emotion recognition. The exclusion of visual cues limits the model's ability to capture non-verbal signals such as facial expressions and gestures, which are often essential for identifying sadness.

Second, the evaluation was conducted using a single benchmark dataset with a fixed training-test split, which may restrict the generalizability of the results across domains, speaker populations, and spontaneous emotional expressions. This design limits the model's exposure to diverse linguistic and acoustic variations, potentially leading to overfitting to dataset-specific patterns.

Third, this study focused exclusively on a CNN-LSTM architecture enhanced with a temporal attention mechanism. While this model showed strong performance, it does not reflect the full range of modern architectures available for multimodal learning. Limiting the exploration to a single architecture constrains the ability to assess whether other models such as transformers or pretrained

multimodal encoders that might capture cross-modal dependencies more effectively or offer better generalization across diverse input conditions.

Finally, the task was limited to binary emotion classification (sad vs. non-sad). While this simplification facilitates model development and allows for clearer performance comparisons, it fails to capture the nuanced spectrum of human emotional expression. In real-world scenarios, emotions are often mixed, subtle, or context-dependent, and cannot be easily reduced to a binary distinction. This limitation restricts the ecological validity and applicability of the model, particularly in domains such as mental health monitoring or empathetic human-computer interaction, where the ability to distinguish between multiple or overlapping emotional states is critical.

## 6.2 Future Work

Future research could address several limitations identified in this study. First, the integration of visual modalities such as facial expressions or gestures could be explored to enhance the richness of emotional cues available to the model. Incorporating visual features may improve the model's ability to detect subtle or ambiguous expressions of sadness, especially in multimodal interaction settings.

Second, to improve the generalizability of results, future work should evaluate model performance on multiple datasets, including those that capture more spontaneous and diverse emotional expressions across different speakers, contexts, and recording conditions. Cross-corpus validation and out-of-distribution testing would offer a more rigorous assessment of model robustness and reduce the risk of overfitting to dataset-specific characteristics.

Third, since the attention mechanism yielded the highest performance in this study, future work may focus on further optimizing attention strategies. Potential avenues include hierarchical or cross-modal attention, modality-specific attention gates, or adaptive attention guided by contextual or speaker-level information. These extensions may enhance the model's ability to selectively emphasize informative temporal and modal segments, especially for minority classes.

Additionally, expanding the emotion classification task from binary to multi-class or even continuous representations (e.g., valence-arousal dimensions) would better reflect the complexity of human emotional expression. This would also broaden the applicability of the model to real-world use cases such as mental health monitoring, conversational agents, and affective computing applications where detecting a wider range of emotional states is critical. Future work may include significance testing and confidence interval estimation to ensure that performance improvements are statistically meaningful rather than due to random variation.

## 6.3 Conclusion

This study set out to evaluate the effectiveness of multimodal deep learning architectures in detecting sadness in native English speech, with a particular focus on comparing unimodal and multimodal input configurations. The central aim was to determine whether integrating textual features alongside audio could yield measurable improvements in classification performance, particularly under conditions of class imbalance.

The experimental results provide preliminary support for this hypothesis. The attention-based CNN-LSTM model, which integrates both audio and text inputs and employs a temporal attention mechanism, achieved the highest performance across all configurations. This represents a relative

improvement of over 30%, suggesting that multimodal fusion may help capture emotional cues that are difficult to detect from audio alone.

These findings point toward possible applications in real-world settings. For instance, in the context of mental health monitoring or early intervention, enhanced detection of sadness in speech could support the development of automated tools that complement human judgment. Similarly, in voice assistants and social robotics, modest improvements in emotion recognition may enable more nuanced and empathetic system responses.

In addition to these performance gains, the study offers a modest methodological contribution by demonstrating that an attention-based CNN-LSTM architecture can be effective in a constrained binary classification setting. While not exhaustive, this approach may be of interest to researchers working with limited data or seeking interpretable, efficient models for affective tasks.

This work has clear limitations, including its focus on a single dataset, binary emotion categories, and one specific model architecture. Nonetheless, the findings may serve as a useful reference point for future studies exploring more advanced architectures, additional modalities, or richer emotional representations. Taken together, the results offer an initial step toward understanding how multimodal deep learning can support the development of more sensitive and context-aware emotion recognition systems.

## 6.4 Declaration of AI Usage

I hereby declare that this Master's thesis is my own original work, and that it has not been submitted, either in whole or in part, for any other degree or professional qualification. All sources used, including printed materials, online content, and other resources, have been appropriately acknowledged and referenced where applicable.

During the preparation of this thesis, I made use of OpenAI's GPT-4o for the following purposes: proofreading the entire document to improve linguistic clarity and ensure alignment with academic writing conventions; generating alternative explanations for technical concepts in Section 3; assisting in the visual enhancement of Figure 1; creating initial templates for code documentation; helping to resolve script execution errors; and summarizing background literature during the preliminary review stage. All AI-assisted content was carefully reviewed, verified, and substantially revised by me to ensure accuracy, originality, and academic integrity.

## Bibliography

- Samuel Albanie, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Emotion recognition in speech using cross-modal transfer in the wild. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 292–301. ACM, 2018. doi: 10.1145/3240508.3240572.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Jianhao Lian et al. A comprehensive review on multimodal emotion recognition: Datasets, methods and challenges. *Information Fusion*, 89:120–142, 2023.
- Seung-Hyun Lim and Su-Youn Lee. Speech emotion recognition using convolutional and recurrent neural networks. In *2016 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*, pages 1–4. IEEE, 2016.
- Taniya Mittal, Utsav Bhattacharya, Rohan Chandra, Aniket Bera, and Dinesh Manocha. M3er: Multiplicative multimodal emotion recognition using facial, textual, and speech cues. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 1359–1367, 2020.
- Maja Pantic and Leon Rothkrantz. Emotion recognition through multiple modalities: Face, body gesture, speech. In Jianhua Tao, Tieniu Tan, and Rosalind Picard, editors, *Affective Information Processing*, pages 113–135. Springer, 2008. doi: 10.1007/978-3-540-85099-1\_8. URL [https://link.springer.com/chapter/10.1007/978-3-540-85099-1\\_8](https://link.springer.com/chapter/10.1007/978-3-540-85099-1_8).
- Soujanya Poria, Erik Cambria, Rajiv Bajpai, and Amir Hussain. A review of affective computing: From unimodal analysis to multimodal fusion. *Information Fusion*, 37:98–125, 2017. doi: 10.1016/j.inffus.2017.02.003. URL <https://doi.org/10.1016/j.inffus.2017.02.003>.
- Takaya Saito and Marc Rehmsmeier. The precision-recall plot is more informative than the roc plot when evaluating binary classifiers on imbalanced datasets. *PLOS ONE*, 2015.
- Rahul Singh, Pooja Sharma, and Raghav Verma. Improving minority class performance in speech emotion recognition using hierarchical learning. In *Proceedings of the 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6259–6263. IEEE, 2021.
- Chamath Siriwardhana, João Reis, Chathura Fernando, and Suranga Nanayakkara. Jointly fine-tuning bert-based representations for multimodal speech emotion recognition. In *Proceedings of Interspeech*, pages 3755–3759, 2020.
- Panagiotis Tzirakis, George Trigeorgis, Mihalis A Nicolaou, Björn Schuller, and Stefanos Zafeiriou. End-to-end multimodal emotion recognition using deep neural networks. In *Proceedings of the 19th ACM International Conference on Multimodal Interaction*, pages 394–401. ACM, 2017.
- Sungjoon Yoon, Seungryong Byun, and Kyomin Jung. Multimodal speech emotion recognition using audio and text. In *2018 IEEE Spoken Language Technology Workshop (SLT)*, pages 112–118. IEEE, 2018. doi: 10.1109/SLT.2018.8639580. URL <https://doi.org/10.1109/SLT.2018.8639580>.