

# Two-Stage Emotion Detection from Multimodal Data

Xiangyi Li

San José State University  
Department of Computer Science

Spring 2024

# Outline

- 1 Introduction
- 2 Background
- 3 Methodology
- 4 Results
- 5 Conclusions

# Outline

- 1 Introduction
- 2 Background
- 3 Methodology
- 4 Results
- 5 Conclusions

- **Emotion detection** is crucial for human-computer interaction
- Enables machines to recognize and respond to human emotional states
- Applications:
  - Mental health monitoring
  - Customer service
  - Human-computer interaction
  - Sentiment analysis
- **Challenge:** Emotions are complex, multidimensional phenomena

# Research Questions

- ① How does a **two-stage approach** (dimensional prediction → category mapping) compare to **direct classification** for emotion recognition?
- ② What is the relative contribution of **text vs. audio modalities** for emotion detection?
- ③ Which **fusion strategies** best integrate multimodal information?
- ④ How do different **transformer architectures** perform for emotion detection tasks?

# Outline

- 1 Introduction
- 2 Background
- 3 Methodology
- 4 Results
- 5 Conclusions

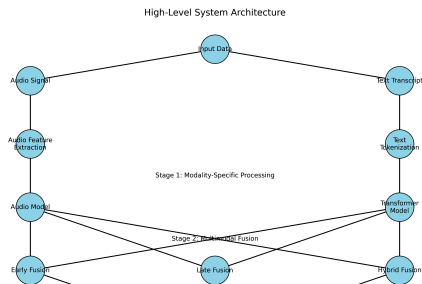
# Dimensional vs. Categorical Emotion Models

## Dimensional Model

- Represents emotions as points in continuous space
- **AVD dimensions:**
  - **A**rousal: energy/intensity
  - **V**alence: positive/negative
  - **D**ominance: control/power
- Captures nuanced emotional states

## Categorical Model

- Discrete emotion labels (anger, joy, sadness, etc.)
- Easier to classify
- More intuitive for humans
- Less granular representation



# Evolution of Emotion Recognition

- **Pre-2012:** Mostly rule-based systems and traditional ML
  - SVM, Decision Trees, Bayesian methods
  - Handcrafted features like lexicons and acoustic parameters
- **Deep Learning Era (2013-2017):**
  - CNNs, RNNs for feature extraction
  - Word embeddings (Word2Vec, GloVe)
- **Transformer Era (2018-Present):**
  - BERT, RoBERTa, XLNet, DeBERTa
  - Attention-based architectures enable better context modeling

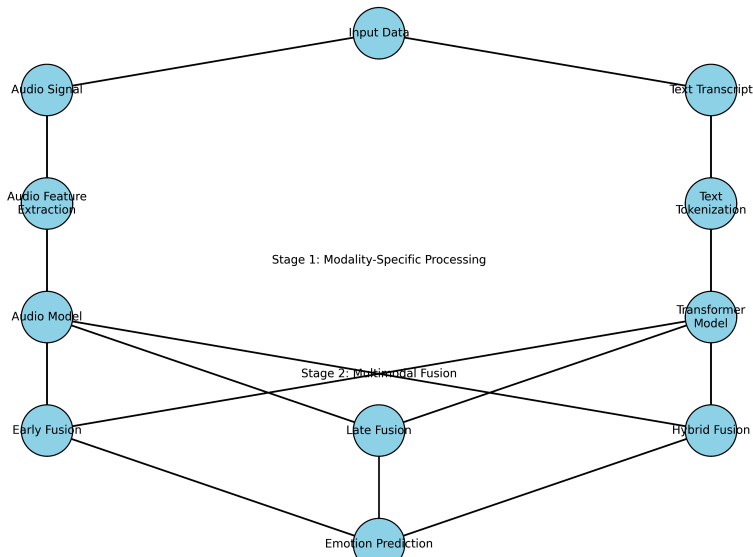


# Outline

- 1 Introduction
- 2 Background
- 3 Methodology**
- 4 Results
- 5 Conclusions

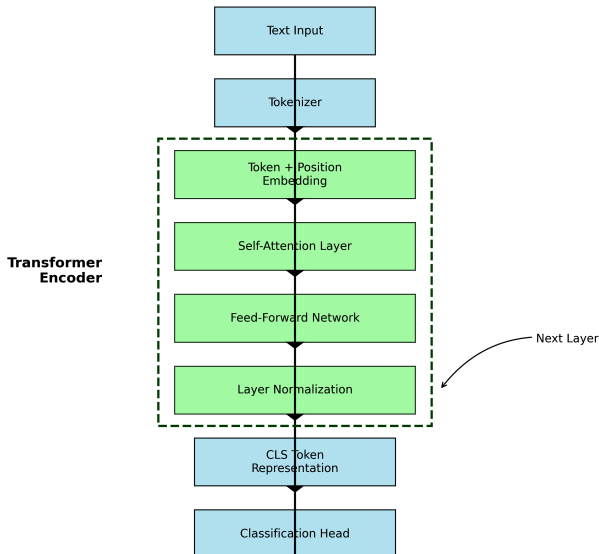
# System Architecture

## High-Level System Architecture



# Text Processing Models

## Text Model Architecture Detail



# Audio Feature Extraction

## Audio Feature Extraction Process



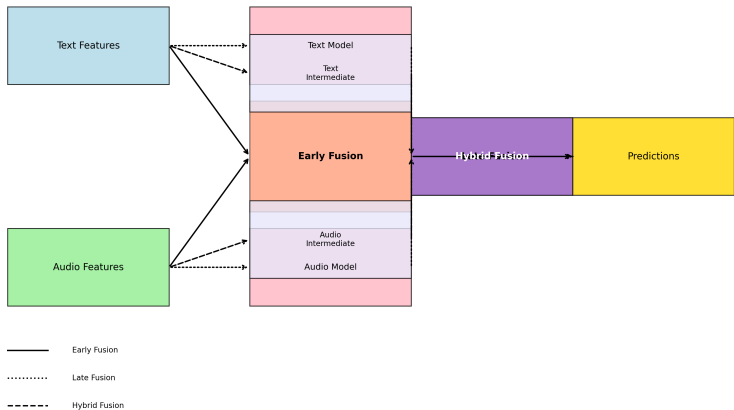
Audio

feature extraction pipeline

- **MFCCs:** Mel-Frequency Cepstral Coefficients (vocal tract characteristics)

# Fusion Strategies

Comparison of Fusion Strategies



## Multimodal fusion approaches

# Experimental Setup

- **Dataset:** IEMOCAP (Interactive Emotional Dyadic Motion Capture)
  - 12 hours of audio-visual data
  - 10 speakers (5 male, 5 female)
  - Both categorical and dimensional annotations
- **Implementation:** PyTorch, Hugging Face Transformers
- **Training Protocol:**
  - AdamW optimizer with linear learning rate schedule
  - Early stopping based on validation loss
  - 5-fold cross-validation
- **Evaluation Metrics:** Accuracy, F1 (Macro/Micro), RMSE, MAE

# Outline

- 1 Introduction
- 2 Background
- 3 Methodology
- 4 Results**
- 5 Conclusions

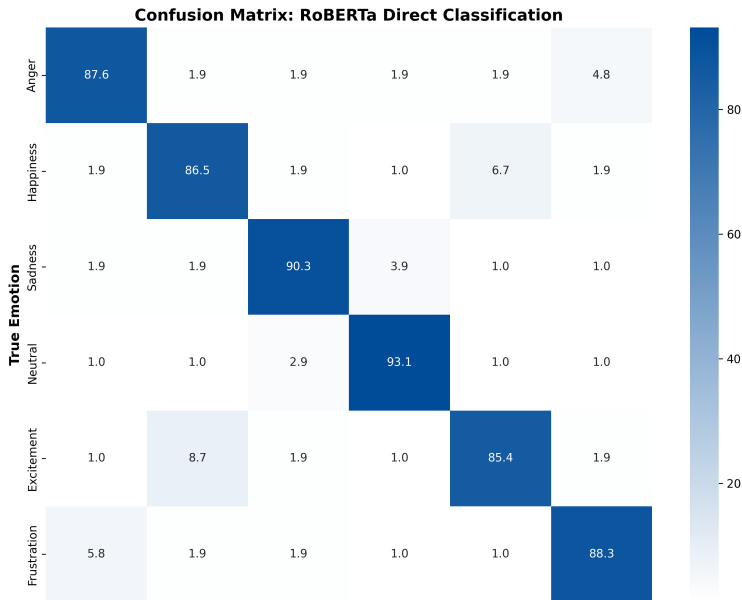
# Dimensional Emotion Prediction (Stage 1)

Model	Modality	Dimension	Test RMSE	MAE
RoBERTa	Text	Valence	0.630	0.500
		Arousal	0.730	0.560
		Dominance	0.680	0.530
CNN+MFCC	Audio	Valence	0.720	0.590
		Arousal	0.650	0.510
		Dominance	0.700	0.560
RoBERTa+MFCC	Multimodal	Valence	0.610	0.490
		Arousal	0.640	0.500
		Dominance	0.660	0.520

- Text models perform better for **Valence** (positive/negative sentiment)
- Audio models perform better for **Arousal** (intensity/energy)
- Multimodal approaches show balanced performance across dimensions



# Categorical Emotion Classification

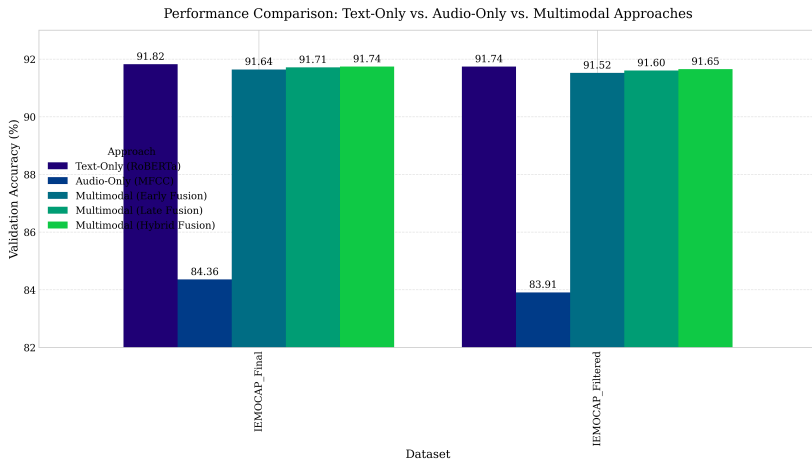


# Two-Stage vs. Direct Classification

Approach	Modality	Test Acc.	Macro F1	M
Direct (RoBERTa)	Text	0.95	0.94	
Two-Stage (RoBERTa)	Text	0.92	0.91	
Direct (CNN+MFCC)	Audio	0.89	0.87	
Two-Stage (CNN+MFCC)	Audio	0.87	0.85	
Direct (RoBERTa+MFCC)	Multimodal	0.94	0.93	
Two-Stage (RoBERTa+MFCC)	Multimodal	0.90	0.89	

- Direct classification consistently outperforms two-stage approach
- But two-stage approach provides richer emotional representation
- Performance gap consistent across modalities (1.5-2.5%)

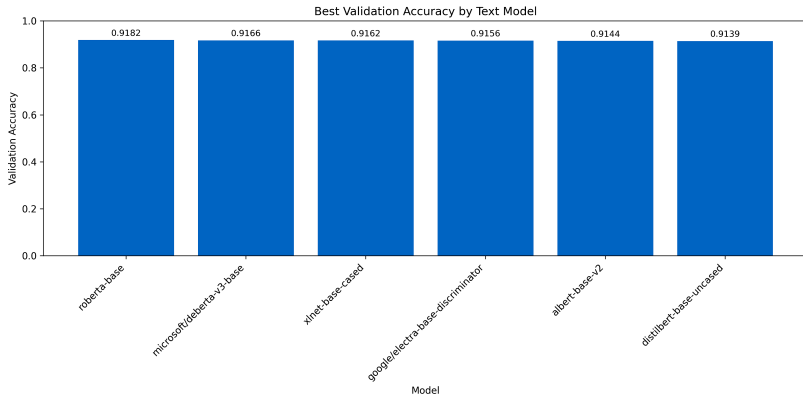
# Modality Importance



## Performance comparison across modalities

- Text-only approaches slightly outperform multimodal approaches
- But gap narrows with optimal fusion strategies

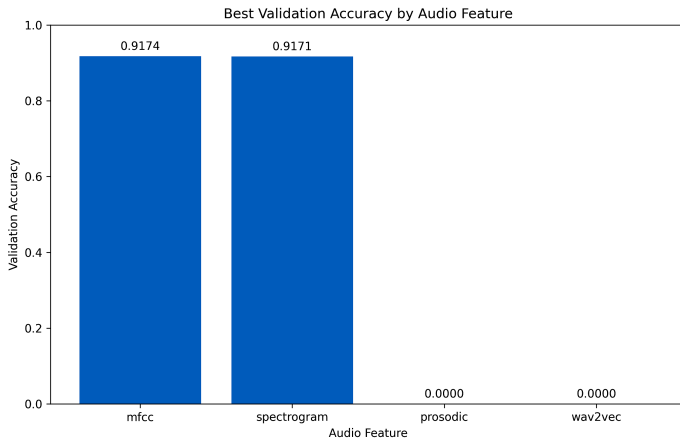
# Transformer Model Comparison



## Performance comparison of transformer models

- RoBERTa consistently outperforms other models
- DeBERTa shows strong performance, particularly for valence
- ALBERT shows lowest performance despite parameter efficiency

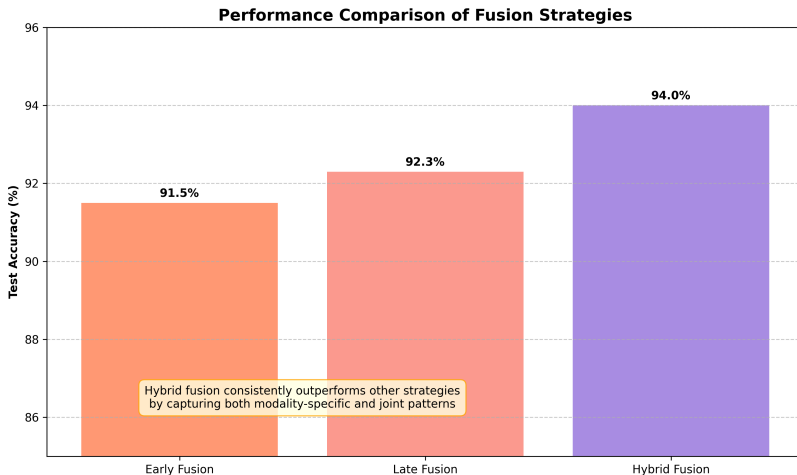
# Audio Feature Effectiveness



## Comparison of audio feature extraction methods

- MFCCs provide the best performance for emotion detection
- Spectrograms capture more temporal information but are noisier

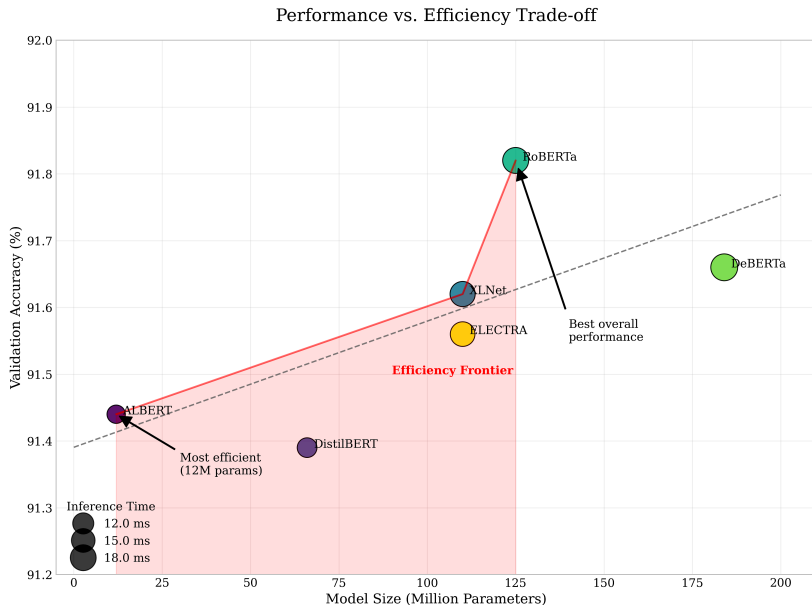
# Fusion Strategy Considerations



## Performance comparison of fusion strategies

- Attention-based fusion provides best overall performance
- Late fusion performs well for categorical classification

# Performance-Efficiency Tradeoffs



# Outline

- 1 Introduction
- 2 Background
- 3 Methodology
- 4 Results
- 5 Conclusions**



# Key Findings

- Direct classification slightly outperforms two-stage approach for categorical emotion recognition
- Text-only approaches slightly outperform multimodal ones, though the gap narrows with optimal fusion
- Textual features better capture valence, while audio features more effectively represent arousal
- RoBERTa consistently outperforms other transformer models
- Attention-based fusion provides the best integration of multimodal information

- **Application-Specific Approach Selection:**

- Direct classification: When accuracy is critical
- Two-stage approach: When continuous emotional representation is valuable

- **Resource Considerations:**

- Text-only approaches offer better efficiency
- ALBERT provides good performance-efficiency tradeoff

- **Modality Selection:**

- Valence-focused applications: Prioritize text
- Arousal-focused applications: Incorporate audio

- Incorporate visual modality (facial expressions, gestures)
- Explore more sophisticated fusion techniques (cross-modal attention)
- Investigate culture-specific emotional expressions
- Develop personalized emotion recognition models
- Explore few-shot and zero-shot learning for emotion recognition
- Evaluate on more diverse datasets across languages and contexts

## Questions?

Contact: [xiangyi.li@sjsu.edu](mailto:xiangyi.li@sjsu.edu)