

# Multimodal Affective Analysis Using Hierarchical Attention Strategy with Word-Level Alignment

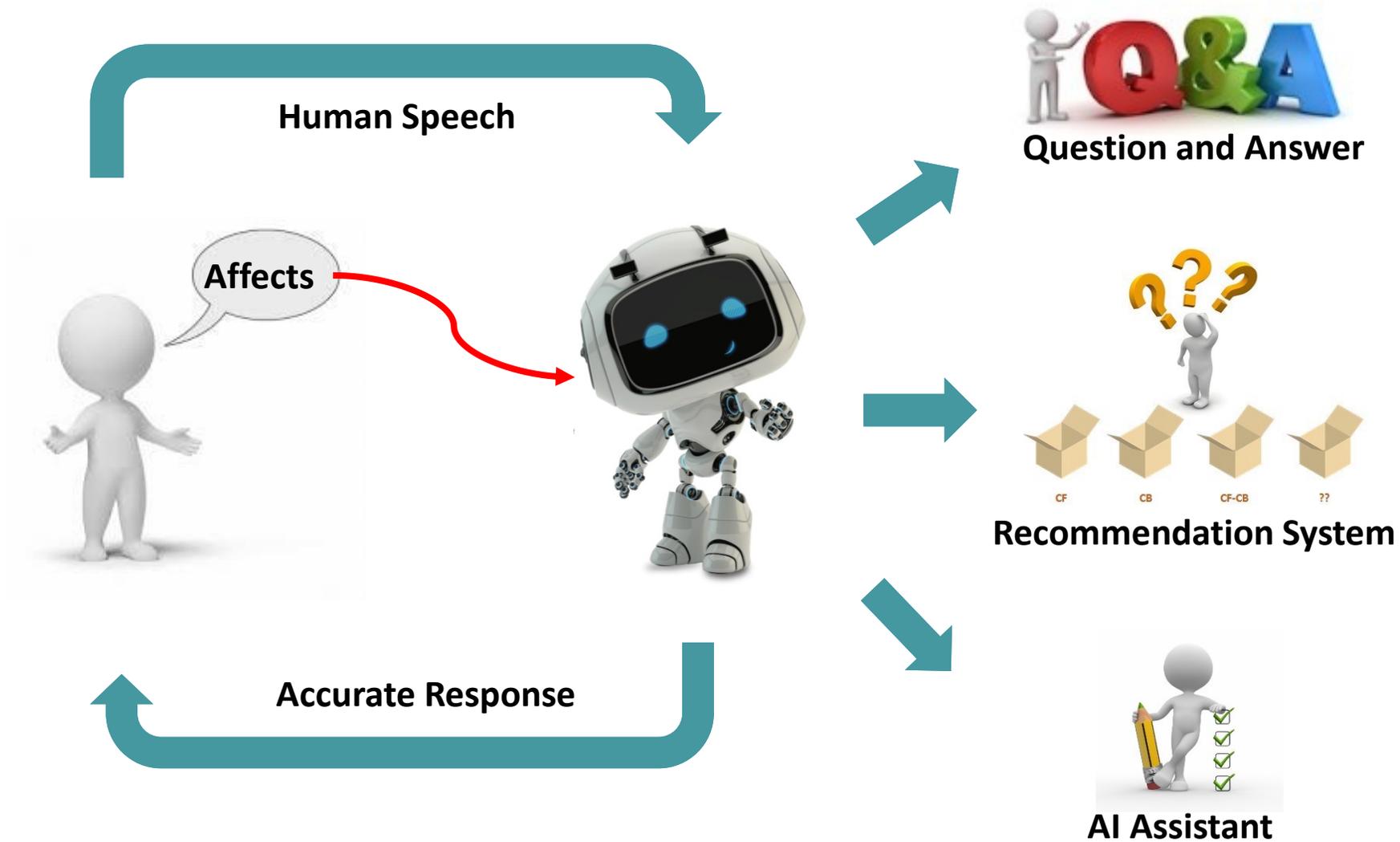
**Yue Gu\***, Kangning Yang, Shiyu Fu, Shuhong Chen, Xinyu Li, Ivan Marsic

Multimedia Image Processing Lab

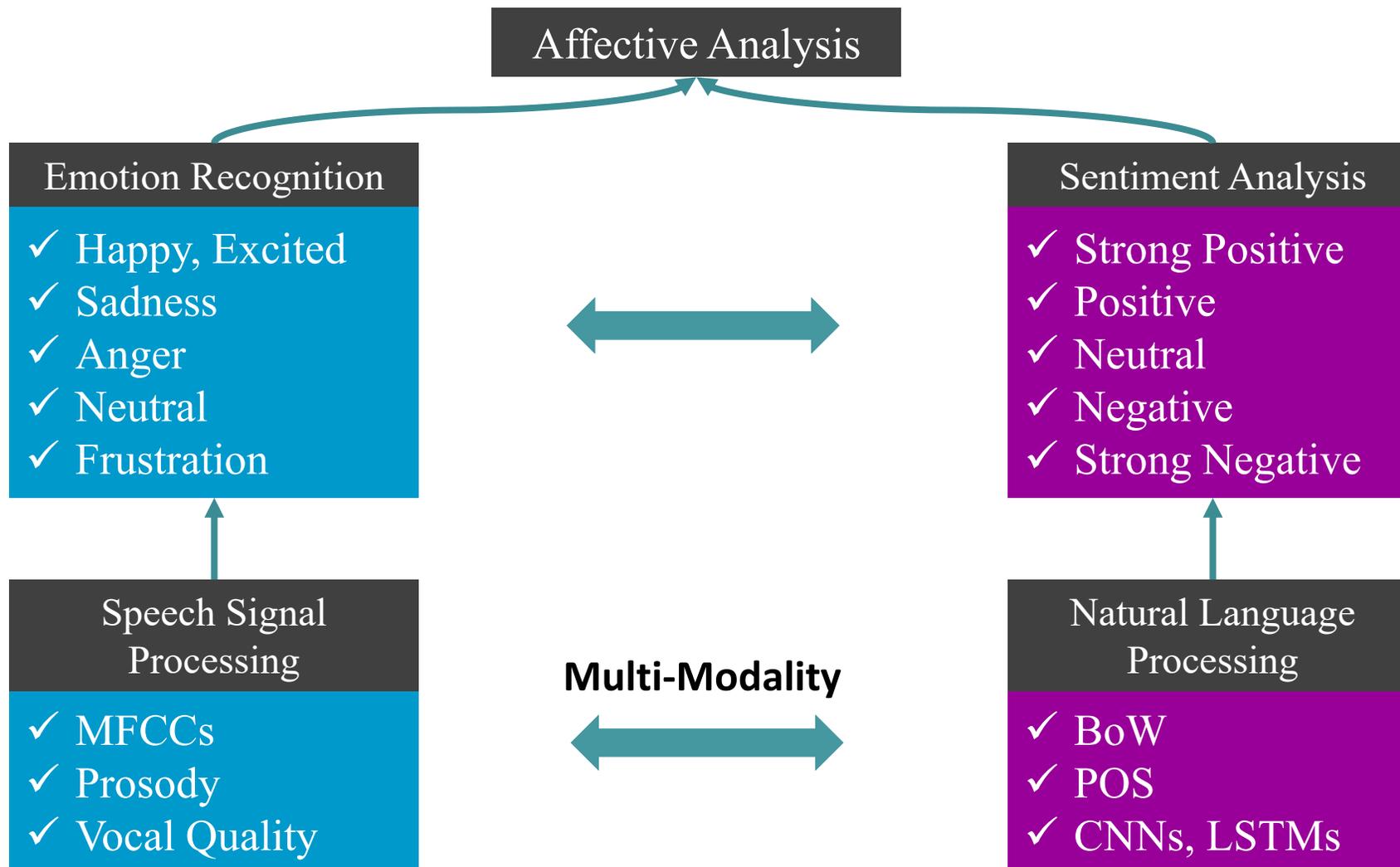
Electrical and Computer Engineering Department

Rutgers, The State University of New Jersey

# Why the affective analysis is necessary?



# Progress of Affective Computing



# Is multi-modality needed?

## ➤ Vocal signal prominence

Oh you don't like that you are west-sider



*Neutral* or *Frustration*

# Is multi-modality needed?

## ➤ Vocal signal prominence

Oh you don't like that you are west-sider



*Happy*

## Is multi-modality needed?

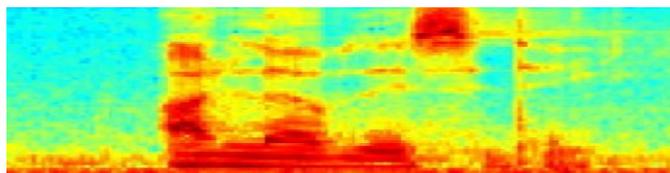
### ➤ Vocal signal prominence

Oh you don't like that you are west-sider

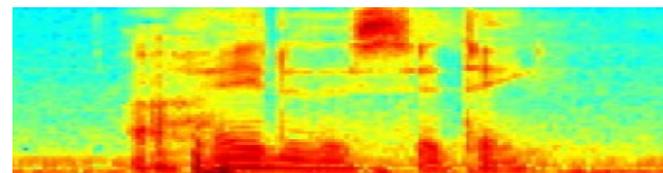


*Happy*

### ➤ Acoustic ambiguity



**"I love this city!"**



**"I hate this city!"**

## Challenges: Feature Extraction

- Gap between features and actual affective states
  - Lack of high-level associations
  - Not all parts contribute equally

## Challenges: Modality Fusion

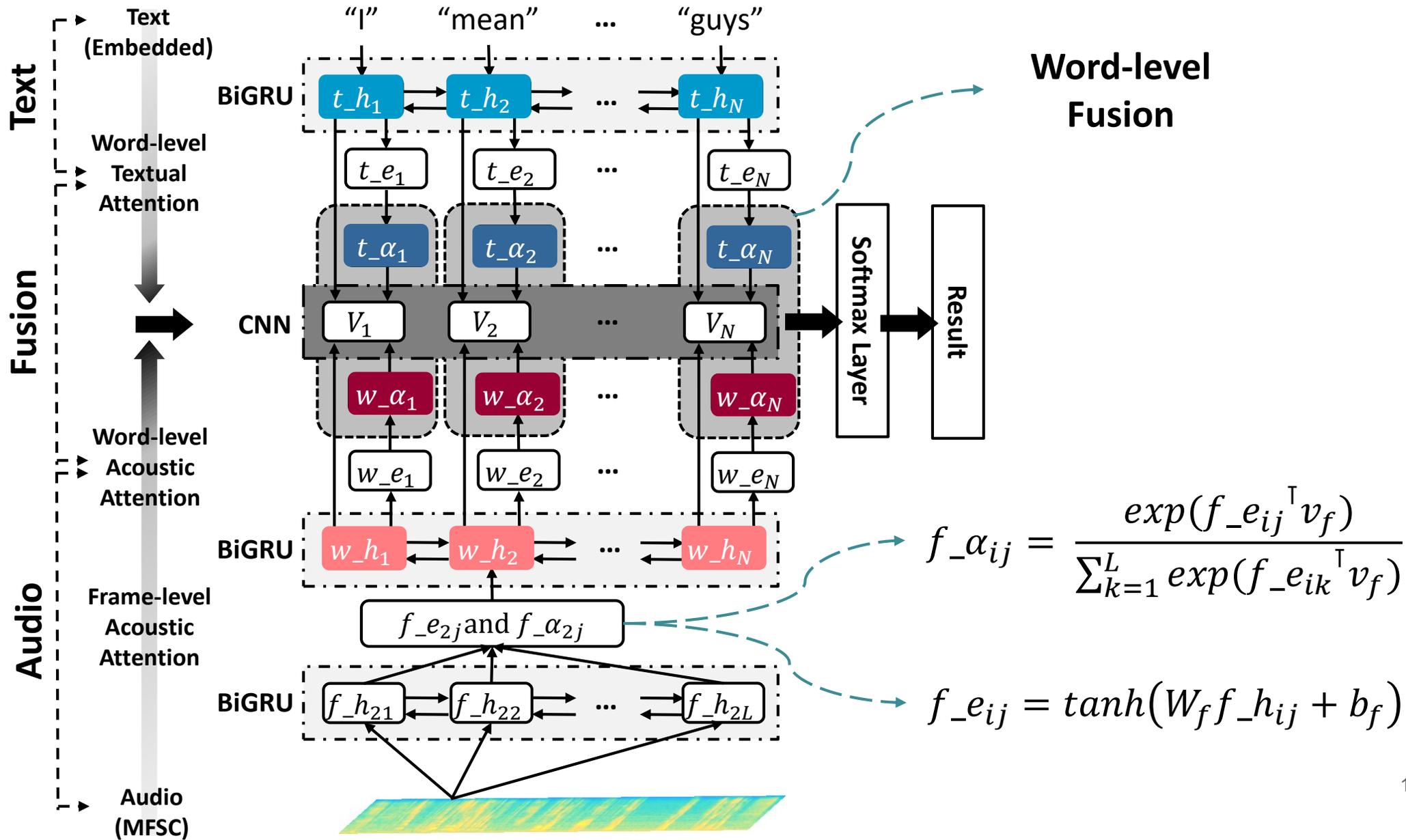
- Decision-level Fusion
  - Lack of mutual association learning
- Feature-level Fusion
  - Fail to learn time-dependent interactions
  - Lack of consistency

## Proposed Solutions

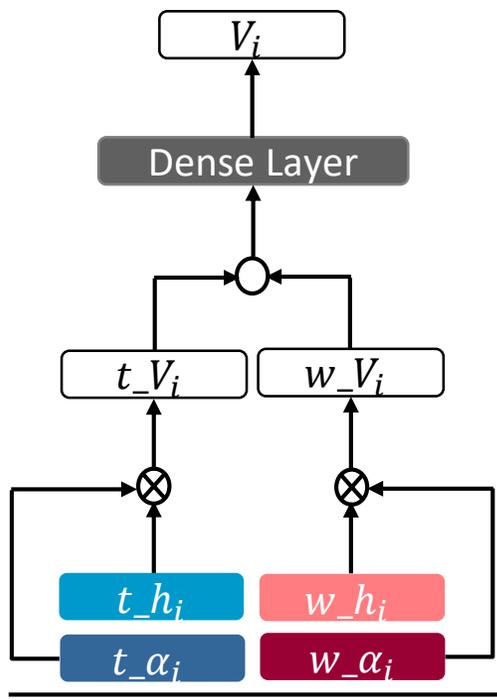
- Feature Extraction
  - Hierarchical attention based bidirectional GRUs
- Modality Fusion
  - Word-level fusion with attention
  - An End-to-End multimodal network

# Data Pre-processing

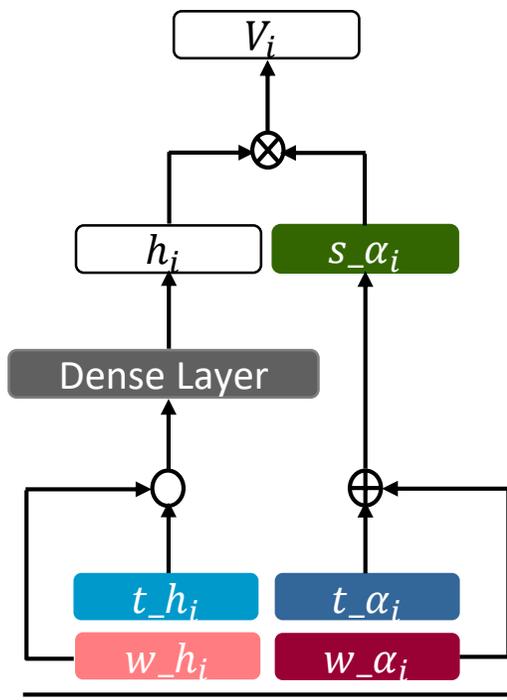
- Text Branch
  - Word Embedding: *word2vec*
- Audio Branch
  - Mel-frequency spectral coefficients (MFSCs)
- Synchronization
  - Word-level forced alignment



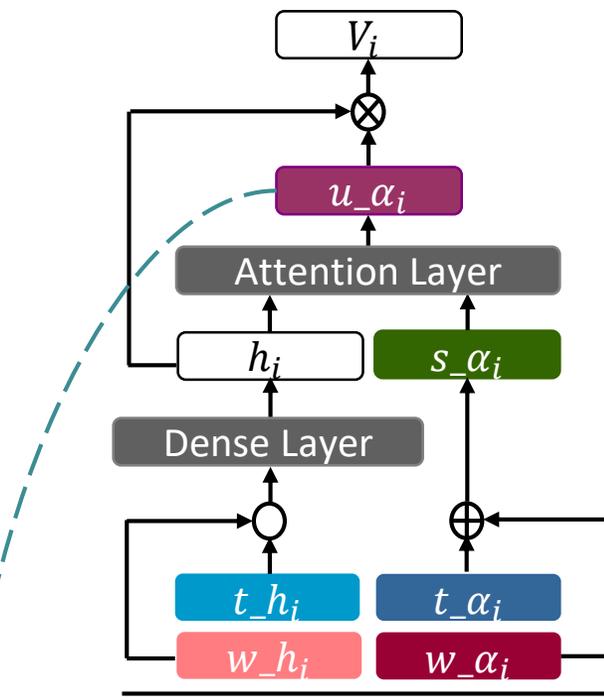
# Word-level Fusion



(a) Horizontal Fusion



(b) Vertical Fusion



(c) Fine-tuning Attention Fusion

- w\_alpha\_i Word-level acoustic attention distribution
- t\_alpha\_i Word-level textual attention distribution
- w\_h\_i Word-level acoustic contextual state
- t\_h\_i Word-level textual contextual state

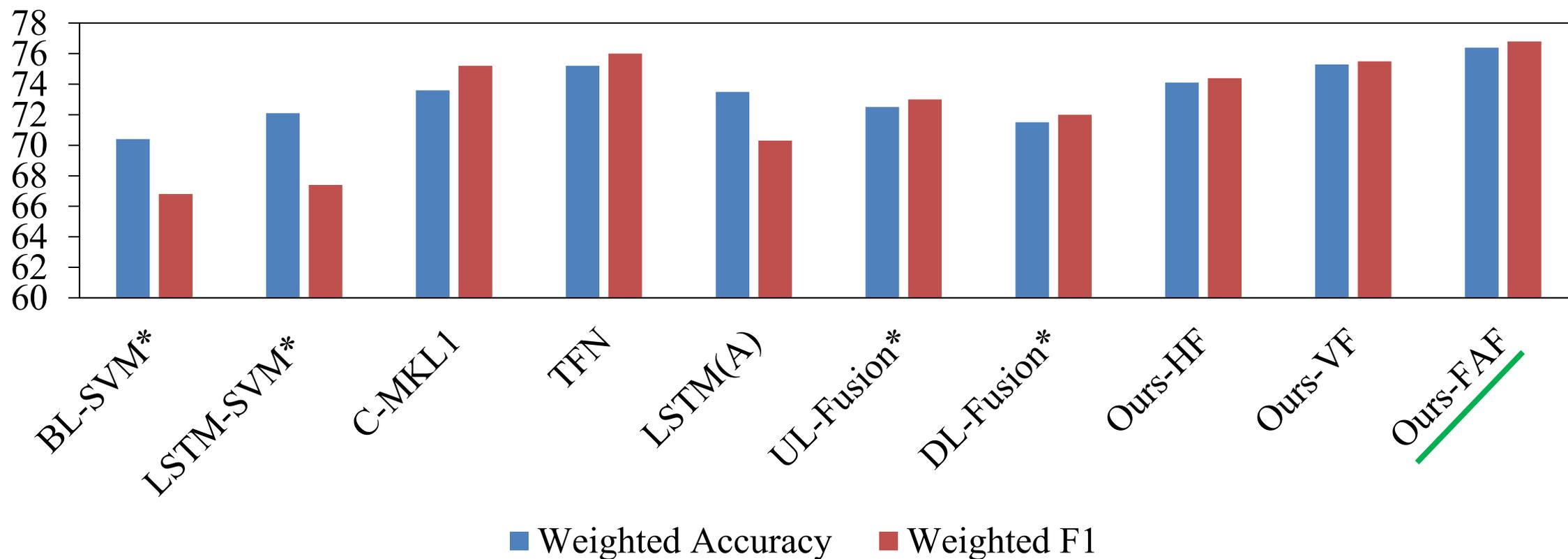
$$u_{\alpha_i} = \frac{\exp(u_{e_i}^T v_u)}{\sum_{k=1}^N \exp(u_{e_k}^T v_u)} + s_{\alpha_i}$$

## Baselines

- Sentiment Analysis
  - BL-SVM, LSTM-SVM
  - C-MKL, TFN, LSTM(A)
- Emotion Recognition
  - SVM Trees, GSV-eVector
  - C-MKL, H-DMS
- Fusion
  - Decision-level, Feature-level (utterance-level)

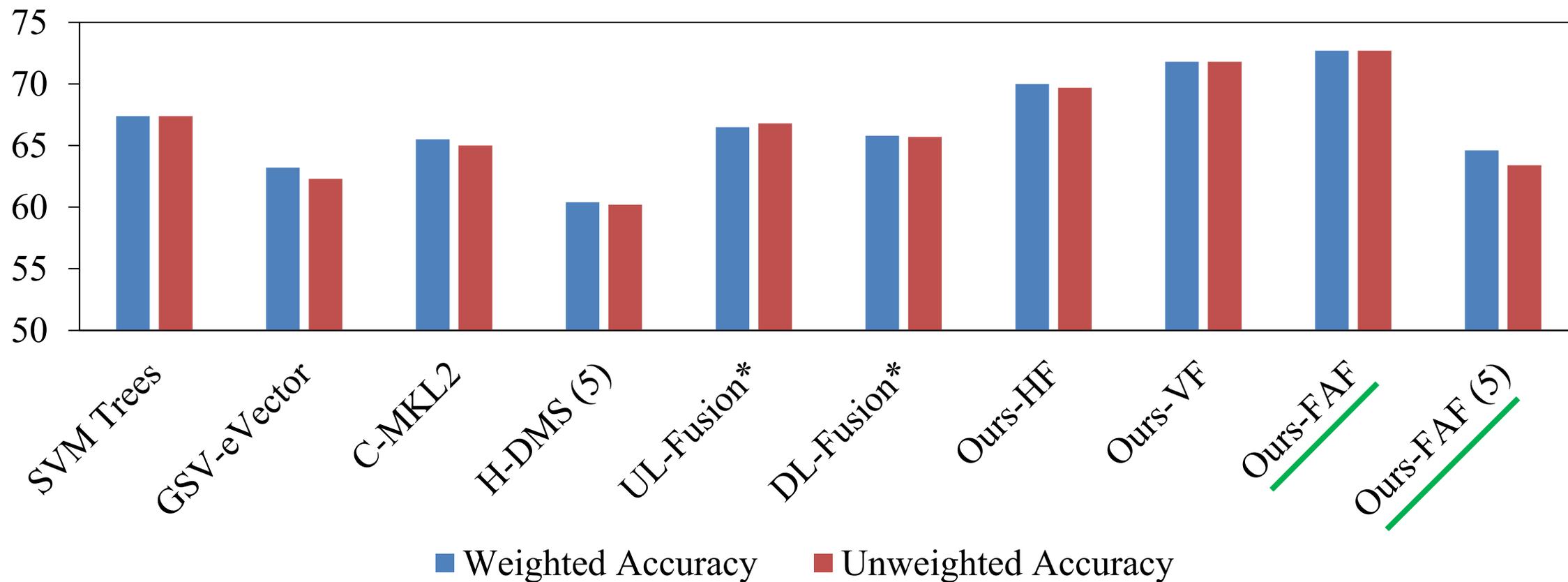
# Sentiment Analysis Result

MOSI



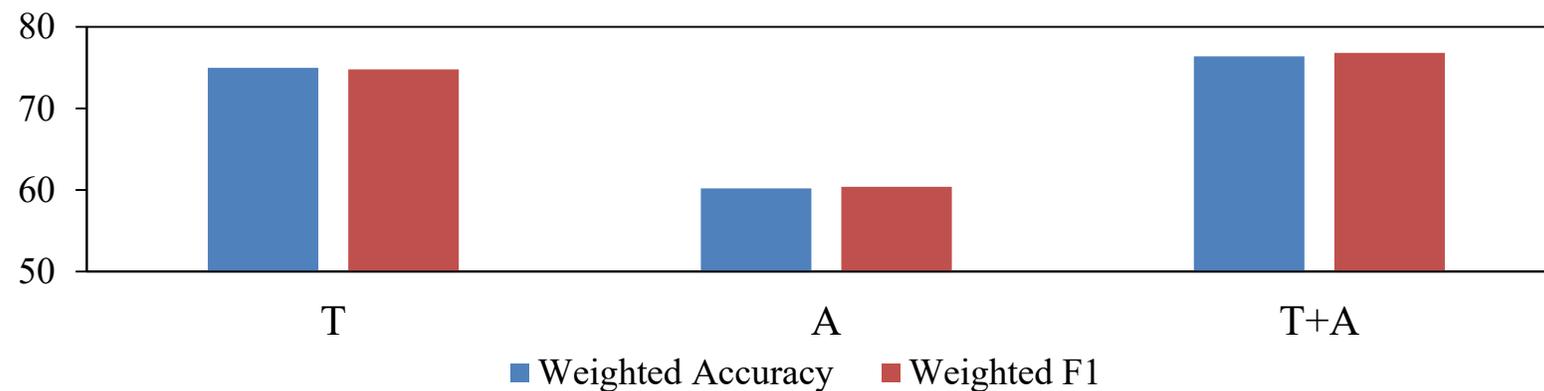
# Emotion Recognition Result

## IEMOCAP

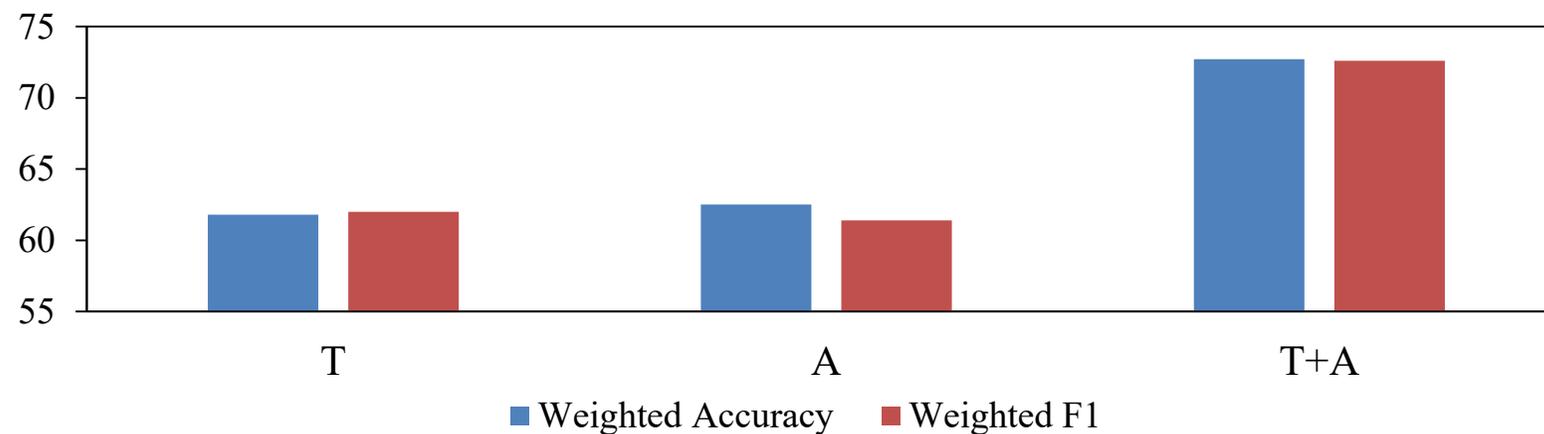


# Multimodal architecture is needed

## MOSI

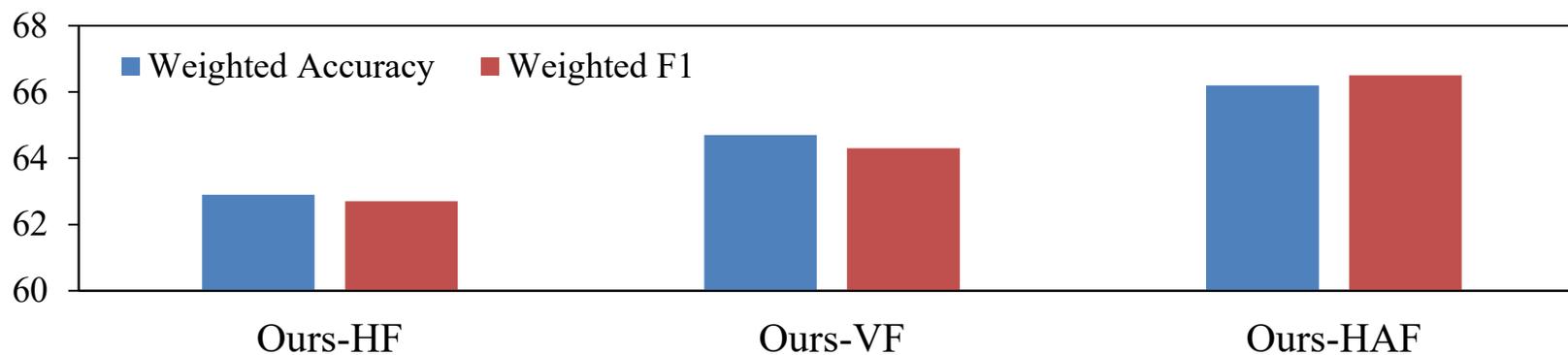


## IEMOCAP

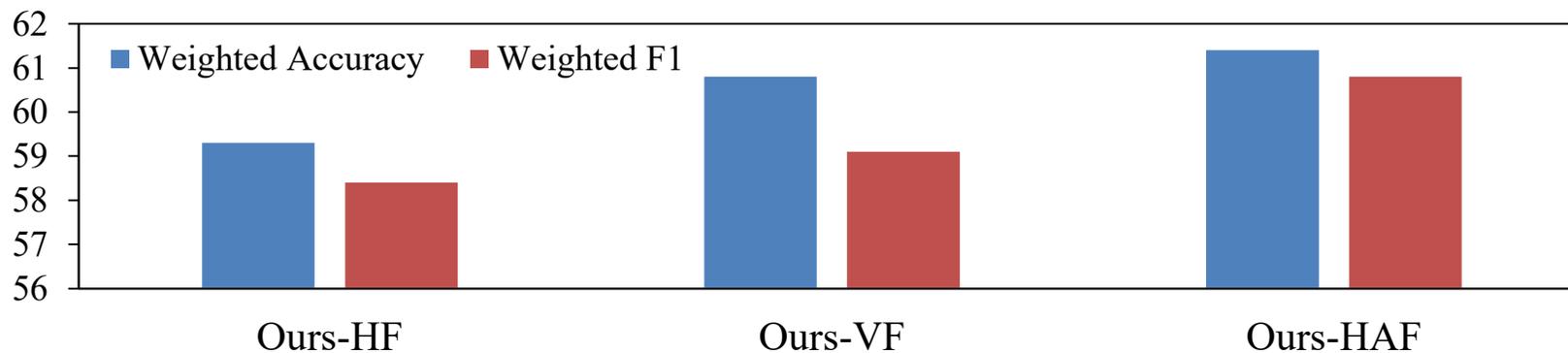


# Generalization

## MOSI to YouTube



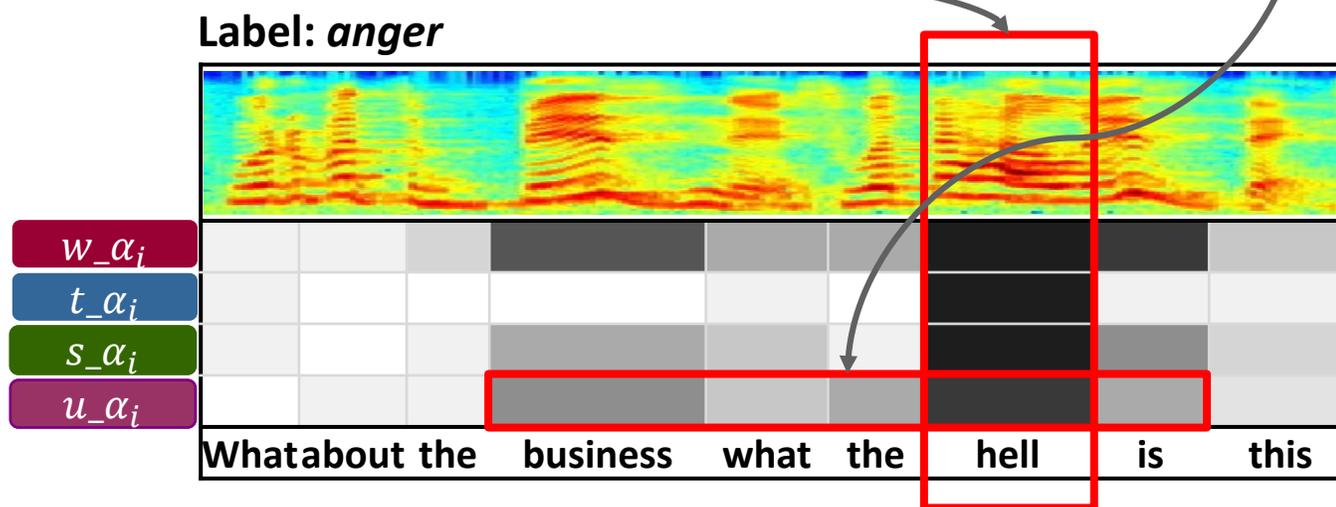
## IEMOCAP to EmotiW



# Attention Visualization

Carry representative information in both text and audio

Successfully combine both textual and acoustic attentions



$w_{\alpha_i}$  Word-level acoustic attention distribution  
 $t_{\alpha_i}$  Word-level textual attention distribution

$s_{\alpha_i}$  Shared attention distribution  
 $u_{\alpha_i}$  Fine-tuning attention distribution



## Summary

- A hierarchical attention based multimodal structure
- The word-level fusion strategies
- Word-level attention visualization

# Thank you !

**Email:** [yg202@scarletmail.rutgers.edu](mailto:yg202@scarletmail.rutgers.edu)

**Homepage:** [www.ieyuegu.com](http://www.ieyuegu.com)