# MULTI-LOSS LEARNING FOR SPEECH EMOTION RECOGNITION WITH ENERGY-ADAPTIVE MIXUP AND FRAME-LEVEL ATTENTION

*Cong Wang*[1,†,‡]    *Yizhong Geng*[1,†,‡]    *Yuhua Wen*[1]    *Qifei Li*[1]    *Yingming Gao*[1]
*Ruimin Wang*[2]    *Chunfeng Wang*[2]    *Hao Li*[2]    *Ya Li*[1,*]    *Wei Chen*[2,*]

[1]Beijing University of Posts and Telecommunications, Beijing, China
[2]Li Auto

## ABSTRACT

Speech emotion recognition (SER) is an important technology in human-computer interaction. However, achieving high performance is challenging due to emotional complexity and scarce annotated data. To tackle these challenges, we propose a multi-loss learning (MLL) framework integrating an energy-adaptive mixup (EAM) method and a frame-level attention module (FLAM). The EAM method leverages SNR-based augmentation to generate diverse speech samples capturing subtle emotional variations. FLAM enhances frame-level feature extraction for multi-frame emotional cues. Our MLL strategy combines Kullback-Leibler divergence, focal, center, and supervised contrastive loss to optimize learning, address class imbalance, and improve feature separability. We evaluate our method on four widely used SER datasets: IEMOCAP, MSP-IMPROV, RAVDESS, and SAVEE. The results demonstrate our method achieves state-of-the-art performance, suggesting its effectiveness and robustness.

***Index Terms***— Speech emotion recognition, multi-loss learning, energy-adaptive mixup, frame-level attention

## 1. INTRODUCTION

Speech emotion recognition (SER) is an important technology in human-computer interaction (HCI), enabling systems to recognize and respond to human emotions, improving user interactions. SER has broad applications in healthcare [1], customer service [2], conversational agents [3], and online education [4]. Despite advances, SER remains challenging due to the complex and subjective nature of human emotions.

Recent research indicates emotions in speech are conveyed not just through linguistic content but also subtle nonverbal cues like tone, rhythm, and energy variations [5, 6]. These characteristics are essential for capturing robust emotional features to enhance SER performance. However, annotating emotional speech data is time-consuming and labor-intensive, leading to limited datasets. This data scarcity restricts the learning capacity of deep learning models, hindering SER system performance.

To address data scarcity, researchers have employed data augmentation techniques, such as noise addition, to enhance SER performance. For instance, An et al. [7] applied additive Gaussian white noise to expand their dataset, achieving notable accuracy improvements. More advanced methods, such as mixup from computer vision (CV), have been adopted for their effectiveness. Kang et al. [8] introduced a label-adaptive mixup for SER, combining speech segments to create mixed-label representations. However, this approach ignores energy-based emotional variations by mixing segments uniformly, a simplification that may overlook critical emotional nuances.

To address these issues, we propose a novel multi-loss learning (MLL) framework. It integrates an **energy-adaptive mixup (EAM)** method to generate diverse speech samples with varied energy levels. Concurrently, a **frame-level attention module (FLAM)** is introduced to refine the extraction of these multi-frame emotional cues. The framework is then optimized by our **MLL** strategy, which combines four key losses: Kullback-Leibler divergence for soft label alignment, focal loss for hard samples, and both center and supervised contrastive (SupCon) losses to improve feature discrimination. We also integrate a context broadcasting (CB) mechanism to maximize feature utilization from FLAM.

We validate our approach on four widely-used SER benchmark datasets: IEMOCAP [9], MSP-IMPROV [10], RAVDESS [11], and SAVEE [12]. Experimental results show our method outperforms state-of-the-art models and exhibits strong generalization across various emotional speech conditions, underscoring its superior performance and robustness. The key contributions of this work are summarized as follows:

- We propose a novel EAM method and design a FLAM to capture robust emotional features. To the best of our knowledge, this is the first approach to incorporate the energy factor into the mixup process for speech data.

- We propose a MLL strategy that, for the first time, integrates the supervised contrastive loss and the center loss for SER. This strategy effectively leverages the latent emotional features, leading to significant performance improvements.

- Extensive experiments conducted on four benchmark datasets demonstrate both the effectiveness and strong generalization capability of our proposed method. Our approach consistently outperforms existing state-of-the-art models across all datasets, highlighting not only its superior performance but also its robustness in diverse emotional speech scenarios.

## 2. METHOD

### 2.1. Model architecture

As shown in Figure 1, our model integrates three core components: an energy-adaptive mixup (EAM) method, a frame-level attention module (FLAM), and a multi-loss learning (MLL) strategy. The EAM method enhances training data by generating mixed speech samples with varied energy levels based on Signal-to-Noise Ratio
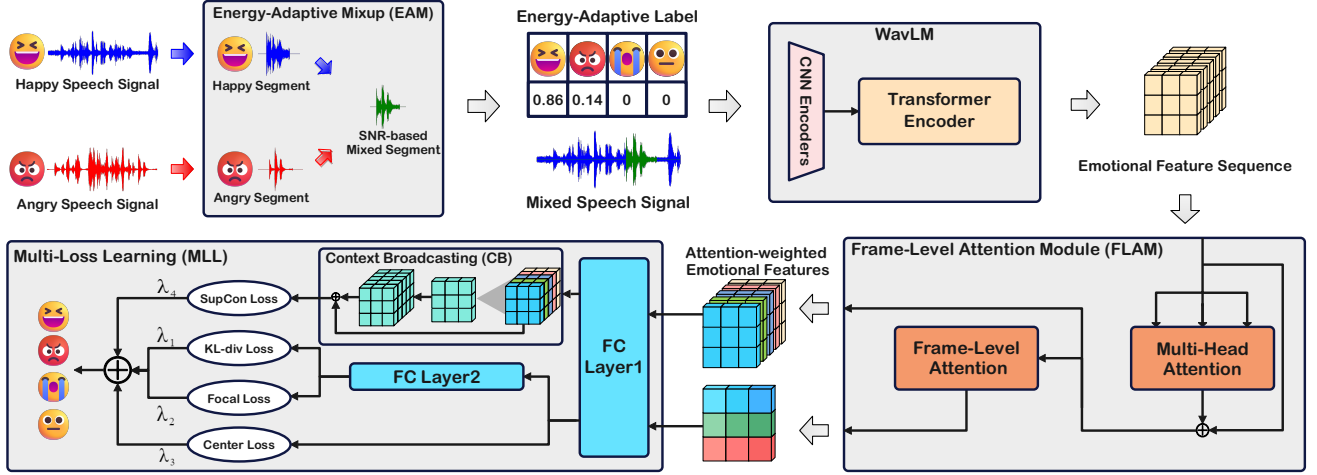
---

**Fig. 1**: The overall model architecture of our proposed SER method.

(SNR), providing richer emotional variations. The FLAM then refines temporal relationships between frames using a multi-head attention mechanism to focus on salient emotional cues. Finally, our MLL strategy optimizes the model by combining multiple specialized loss functions (KL-div, focal, center, and SupCon) to handle label distributions, difficult samples, and improve feature separability.

### 2.2. Energy-adaptive mixup method

Building upon the length-adaptive mixup (LAM) method [8], which overlooks energy factors, we propose our EAM method. In contrast to LAM's length-based label weighting, EAM incorporates energy. The process begins by selecting a random mix length $l_{\text{mix}}$ and starting positions $x_i, x_j$ to extract speech segments $\mathbf{x}'_i$ and $\mathbf{x}'_j$. The energy $P'_j$ of segment $\mathbf{x}'_j$ is then adjusted based on a random SNR(dB) value to generate a new energy $P''_j$, as defined by:

$$\text{SNR(dB)} = 10 \log_{10}\left(\frac{P'_i}{P''_j}\right), \tag{1}$$

From this, the scaling factor and the energy-adjusted speech segment $\mathbf{x}''_j$ are computed:

$$scale = \sqrt{\frac{P'_i}{10^{\text{SNR(dB)}/10} \times P'_j}}, \tag{2}$$

$$\mathbf{x}''_j = scale \times \mathbf{x}'_j. \tag{3}$$

The adjusted segment $\mathbf{x}''_j$ is then mixed with the original speech $\mathbf{x}_i$ by overwriting the corresponding segment:

$$\mathbf{x}_{\text{eng}}[x_i : x_i + l_{\text{mix}}] = \mathbf{x}'_i + \mathbf{x}''_j. \tag{4}$$

This produces an energy-adaptive label that incorporates both length and energy ratios:

$$\mathbf{y}_{\text{eng}} = \left[1 - \frac{l_{\text{mix}}}{l_i}\frac{P''_j}{P'_i + P''_j}, \frac{l_{\text{mix}}}{l_i}\frac{P''_j}{P'_i + P''_j}, \dots, 0\right]. \tag{5}$$

The resulting mixed signal $\mathbf{x}_{\text{eng}}$ is then fed into a pre-trained WavLM Large model [13] to extract an emotional feature sequence $X \in$ $\mathbb{R}^{T \times D}$. By creating richer samples and a more representative label distribution, EAM better reflects the correlation between energy and emotion [14].

### 2.3. Frame-level attention module

Our FLAM enhances inter-frame relationships to capture subtle temporal dependencies. The input feature sequence $X \in \mathbb{R}^{T \times D}$ is first processed by a 16-head Multi-Head Self-Attention (MSA) module with a residual connection to produce an enhanced feature sequence $X'$:

$$X' = X + MultiHead(X). \tag{6}$$

This sequence $X'$ is then aggregated into a single feature vector $f \in \mathbb{R}^D$ using frame-level attention weights learned from a fully connected (FC) layer, a process more effective than traditional pooling:ß

$$f = X'^T \cdot FC(X') \quad \text{where} \quad FC(X') \in \mathbb{R}^{T \times 1}. \tag{7}$$

This module allows the model to aggregate multi-frame features by weighting their relative importance, thus extracting more robust features.

### 2.4. Multi-loss learning strategy

We employ a weighted MLL strategy to optimize the extracted features. Since the EAM label $\mathbf{y}_{\text{eng}}$ is a soft probability distribution, we use KL-divergence loss to measure the difference between the predicted and target distributions:

$$\mathcal{L}_{\text{KL}} = \sum_{i=1}^{|E|} y_i \log\left(\frac{y_i}{\hat{y}_i}\right). \tag{8}$$

To focus on harder-to-classify samples, we apply focal loss [15]:

$$\mathcal{L}_{\text{Focal}} = -\sum_{i=1}^{|E|}(1 - \hat{y}_i)^\gamma \log(\hat{y}_i) y_i. \tag{9}$$

To improve feature discrimination, we project features to a lower-dimensional space. The center loss [16] minimizes intra-class vari-

ance by pulling features towards their corresponding class centers:

$$\mathcal{L}_{\text{Center}} = \frac{1}{B} \sum_{i=1}^{B} \|f_{\text{low}_i} - c_i\|^2. \tag{10}$$

We also adapt a context broadcasting (CB) [17] mechanism to encourage sparse feature interactions before applying SupCon loss:

$$CB(X'_{\text{low}_i}) = \frac{1}{2} \left( X'_{\text{low}_i} + \frac{1}{T} \sum_{m=1}^{T} X'_{\text{low}_m} \right). \tag{11}$$

Finally, supervised contrastive (SupCon) loss [18] maximizes interclass distance while minimizing intra-class similarity:

$$\mathcal{L}_{\text{SupCon}} = \frac{\sum_{i \in I} \frac{-1}{|P(i)|} \sum_{j \in P(i)} \log \frac{\exp(X'_{\text{low}_i} \cdot X'_{\text{low}_j}/\tau)}{\sum_{k \in A(i)} \exp(X'_{\text{low}_i} \cdot X'_{\text{low}_k}/\tau)}}{B \times T}. \tag{12}$$

Our final objective is a weighted combination of these four loss functions:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{KL}} + \lambda_2 \mathcal{L}_{\text{Focal}} + \lambda_3 \mathcal{L}_{\text{Center}} + \lambda_4 \mathcal{L}_{\text{SupCon}}. \tag{13}$$

## 3. EXPERIMENTS AND RESULTS

### 3.1. Datasets

We evaluate our method on four widely-used SER benchmark datasets: IEMOCAP [9], MSP-IMPROV [10], RAVDESS [11], and SAVEE [12]. These datasets provide a robust evaluation by including a mix of acted and spontaneous emotions, diverse speakers, and varying recording conditions.

**IEMOCAP** [9] contains roughly 12 hours of spontaneous and acted conversational data from 10 speakers. We use 5,531 utterances across four classes (happy, angry, sad, neutral) and perform 5-fold session-independent cross-validation.

**MSP-IMPROV** [10] is a conversational corpus containing 8,438 clips of spontaneous emotions from 12 actors in 6 sessions. We use samples from four basic emotions (happy, angry, sad, neutral) and conduct 6-fold session-independent cross-validation.

**RAVDESS** [11] is a dataset of acted emotions from 24 professional actors. We use the speech portions from 1,440 clips across eight emotions (happy, angry, sad, neutral, fearful, disgust, surprised, calm) and conduct 6-fold subject-independent cross-validation.

**SAVEE** [12] contains 480 utterances from 4 male speakers expressing seven emotions (happy, angry, sad, neutral, fearful, disgust, surprised). We perform 4-fold speaker-independent cross-validation.

### 3.2. Experiment setup

The performance of our method is evaluated under the speaker-independent setting. To ensure a fair comparison with previous SOTA methods, we adopt the widely used unweighted accuracy (UA) and weighted accuracy (WA) metrics. In the ablation studies, we systematically removed key components (EAM, FLAM) and individual loss functions (focal, SupCon) to assess their contributions to the overall performance. Each ablation experiment follows the same training and testing settings as the full mode. The dimensionality of the low-dimensional space used for the central loss and the SupCon loss is set to 64. The batch size is set to 16 for all experiments. We use Adam as the optimizer, and the learning rates are initialized to $1 \times 10^{-4}$ for the model and $5 \times 10^{-3}$ for updating the centers, with all learning rates of each epoch decreasing to 7/8 of its previous rate until the $20^{th}$ epoch. All the experiments are conducted on a NVIDIA RTX3090.

**Table 1**: Comparison results on IEMOCAP, MSP-IMPROV, and RAVDESS datasets (A: Audio-only, M: Multi-modal).

| Dataset | Methods | Modality | WA(%) | UA(%) |
|---|---|---|---|---|
| **IEMOCAP** | Tang *et al.* [19] | A | 71.64 | 72.72 |
| | Sun *et al.* [20] | A | 72.86 | 72.85 |
| | Wang *et al.* [21] | A | 73.37 | 74.18 |
| | He *et al.* [22] | A | 73.80 | 74.25 |
| | Gao *et al.* [23] | A | 74.94 | 76.10 |
| | Kang *et al.* [8] | A | 75.37 | 76.04 |
| | He *et al.* [24] | M | 74.50 | 75.00 |
| | Wang *et al.* [25] | M | 75.20 | 76.40 |
| | **Ours** | A | **78.47** | **79.14** |
| **MSP-IMPROV** | Guo *et al.* [26] | A | 46.20 | 44.70 |
| | Nediyanchath *et al.* [27] | A | 47.30 | 46.10 |
| | Xu *et al.* [28] | A | 47.90 | 45.80 |
| | Cao *et al.* [29] | A | 50.70 | 49.90 |
| | Liu *et al.* [30] | A | 51.51 | 41.56 |
| | Liu *et al.* [31] | A | 55.80 | 55.30 |
| | **Ours** | A | **58.55** | **58.34** |
| **RAVDESS** | Baevski *et al.* [32] | A | 74.38 | 73.44 |
| | Sun *et al.* [33] | A | 72.29 | 70.38 |
| | Chen *et al.* [13] | A | 75.36 | 75.28 |
| | Yu *et al.* [34] | A | 81.86 | 82.75 |
| | Chumachenko *et al.* [35] | M | 79.20 | - |
| | Sadok *et al.* [36] | M | 84.80 | - |
| | Sun *et al.* [33] | M | 87.99 | 87.96 |
| | **Ours** | A | **93.40** | **92.28** |

### 3.3. Results and discussions

Our proposed method consistently outperforms existing state-of-the-art (SOTA) approaches, demonstrating strong performance and generalization across all four benchmark datasets. This success underscores the model's ability to handle both spontaneous and acted emotional speech effectively.

As shown in Table 1, our model achieves a Weighted Accuracy (WA) and Unweighted Accuracy (UA) of 78.47%/79.14% on IEMO-CAP and 58.55%/58.34% on MSP-IMPROV, successfully navigating the challenges of conversational and spontaneous speech. The performance on datasets with acted emotions is also remarkable, reaching 93.40%/92.28% on RAVDESS. On the SAVEE dataset (Table 2), our model achieves an average UA of 72.3%, outperforming the previous SOTA and showing significant gains for challenging speakers.

This superior performance is attributed to the synergy between our core components. The Energy-Adaptive Mixup (EAM) provides diverse training data, the Frame-Level Attention Module (FLAM) captures subtle temporal cues, and the Multi-Loss Learning (MLL) strategy ensures robust optimization and feature discrimination. These results confirm our model's effectiveness and robustness across diverse emotional scenarios without relying on dataset-specific optimizations or multi-modal inputs.

### 3.4. Ablation study

To validate the contribution of each component in our proposed method, we conduct a systematic ablation study on the IEMOCAP dataset, with results presented in Table 3.
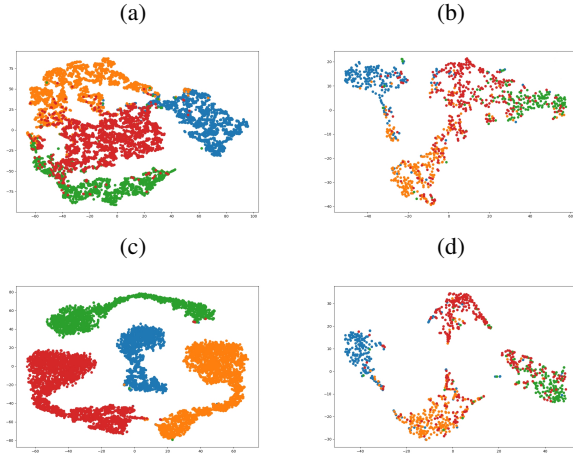
The study first confirms the individual effectiveness of our

**Table 2**: Comparison results on SAVEE in UA(%).

| Methods | DC | JE | JK | KL | Mean |
|---|---|---|---|---|---|
| Chen *et al.* [37] | 81.6 | 83.3 | **69.9** | 49.7 | 71.1 |
| **Ours** | **82.3** | **87.0** | 66.9 | **53.0** | **72.3** |
| Human [37] | 73.7 | 67.7 | 71.2 | 53.2 | 66.5 |

**Table 3**: Ablation study of the proposed method on IEMOCAP, including pre-trained model, mixup method, feature aggregation method, and loss function. We use Kang et al.'s method [8] as the baseline model, as shown in the first row of the table.

| Hubert Large | WavLM Large | Mixup | | Feature Aggregation | | | Loss Function | | | | WA(%) | UA(%) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | LAM | EAM | MaxPool | MeanPool | FLAM | KL-div | Center | Focal | SupCon | | |
| ✓ | – | ✓ | – | – | – | – | ✓ | ✓ | – | – | 75.37 | 76.04 |
| ✓ | – | – | – | – | – | ✓ | ✓ | ✓ | – | – | 75.83 | 76.45 |
| ✓ | – | – | ✓ | – | – | – | ✓ | ✓ | – | – | 76.18 | 76.69 |
| ✓ | – | ✓ | – | – | – | ✓ | ✓ | ✓ | – | – | 76.22 | 76.84 |
| ✓ | – | – | ✓ | – | – | ✓ | ✓ | ✓ | – | – | 76.31 | 76.90 |
| – | ✓ | ✓ | – | – | – | – | ✓ | ✓ | – | – | 76.98 | 77.14 |
| – | ✓ | – | – | – | – | ✓ | ✓ | ✓ | – | – | 77.12 | 77.58 |
| – | ✓ | – | ✓ | – | – | – | ✓ | ✓ | – | – | 77.26 | 77.71 |
| – | ✓ | – | ✓ | ✓ | – | – | ✓ | ✓ | – | – | 77.32 | 77.74 |
| – | ✓ | – | ✓ | – | ✓ | – | ✓ | ✓ | – | – | 77.47 | 77.95 |
| – | ✓ | ✓ | – | – | – | ✓ | ✓ | ✓ | – | – | 77.58 | 78.02 |
| – | ✓ | – | ✓ | – | – | ✓ | ✓ | ✓ | – | – | 77.63 | 78.10 |
| – | ✓ | – | ✓ | – | – | ✓ | ✓ | ✓ | ✓ | – | 77.96 | 78.56 |
| – | ✓ | – | ✓ | – | – | ✓ | ✓ | ✓ | – | ✓ | 78.23 | 78.83 |
| – | ✓ | – | ✓ | – | – | ✓ | ✓ | ✓ | ✓ | ✓ | **78.47** | **79.14** |

(a)       (b)



(c)       (d)

**Fig. 2**: t-SNE visualizations of feature distributions on IEMOCAP. (a) Training set before MLL; (b) Test set before MLL; (c) Training set after MLL; (d) Test set after MLL. Colors: Blue–Angry, Orange–Happy, Green–Sad, Red–Neutral. The feature clusters are visibly more distinct after our MLL strategy.

Energy-Adaptive Mixup (EAM) and Frame-Level Attention Module (FLAM). Results indicate that EAM, which generates energy-diverse samples, provides more meaningful features for aggregation by FLAM than the baseline Length-Adaptive Mixup (LAM). Furthermore, our FLAM significantly outperforms traditional pooling methods like MaxPool and MeanPool by effectively weighting the importance of multi-frame emotional features, leading to a more robust feature representation.

Next, we analyze the components of our Multi-Loss Learning (MLL) strategy. The inclusion of focal loss improves performance by up-weighting hard-to-classify samples. The supervised contrastive (SupCon) loss further enhances feature discrimination by increasing inter-class distance while minimizing intra-class variance. The effectiveness of our MLL strategy in creating more separable feature clusters is visually confirmed by t-SNE visualizations, as shown in Figure 2. Collectively, these results demonstrate the crucial role each component plays in significantly improving SER performance.

## 4. CONCLUSION

In this paper, we introduced a novel Multi-Loss Learning (MLL) framework for Speech Emotion Recognition (SER), integrating an Energy-Adaptive Mixup (EAM) method to generate diverse, energy-varied speech samples and a Frame-Level Attention Module (FLAM) to refine emotional feature extraction. Our MLL strategy effectively addresses class imbalance and improves feature separability. Ablation studies validated the significant contribution of each component, confirming that our integrated framework achieves state-of-the-art (SOTA) performance and robust generalization across multiple benchmark datasets. Future work will focus on validating the model on more diverse cross-linguistic datasets, incorporating multi-modal features (e.g., visual and textual data), and enhancing the mixup method's adaptability through techniques like meta-learning.

# 5. REFERENCES

[1] Agustinus Bimo Gumelar et al., "BiLSTM-CNN hyperparameter optimization for speech emotion and stress recognition," in *2021 International Electronics Symposium (IES)*, 2021, pp. 156–161.

[2] Gudmalwar Ashishkumar Prabhakar et al., "Multichannel CNN-BLSTM architecture for speech emotion recognition system by fusion of magnitude and phase spectral features using DCCA for consumer applications," *IEEE Transactions on Consumer Electronics*, vol. 69, no. 2, pp. 226–235, 2023.

[3] Jiaxiong Hu et al., "The acoustically emotion-aware conversational agent with speech emotion recognition and empathetic responses," *IEEE Transactions on Affective Computing*, vol. 14, no. 1, pp. 17–30, 2023.

[4] Aparna Vyakaranam, Bavani Ramayah, and Tomas Maul, "Preliminary study: Speech emotion recognition in online teaching from the perspective of educators especially late deafened," in *2024 2nd International Conference on Software Engineering and Information Technology (ICoSEIT)*, 2024, pp. 216–221.

[5] Mohammed Jawad Al-Dujaili and Abbas Ebrahimi-Moghadam, "Speech emotion recognition: a comprehensive survey," *Wireless Personal Communications*, vol. 129, no. 4, pp. 2525–2561, 2023.

[6] Kamaldeep Kaur and Parminder Singh, "Trends in speech emotion recognition: a comprehensive survey," *Multimedia Tools and Applications*, pp. 1–45, 2023.

[7] Xu Dong An and Zhou Ruan, "Speech emotion recognition algorithm based on deep learning algorithm fusion of temporal and spatial features," in *Journal of Physics: Conference Series*, 2021, vol. 1861, p. 012064.

[8] Lei Kang, Lichao Zhang, and Dazhi Jiang, "Learning robust self-attention features for speech emotion recognition with label-adaptive mixup," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.

[9] Carlos Busso et al., "IEMOCAP: Interactive emotional dyadic motion capture database," *Language Resources and Evaluation*, vol. 42, pp. 335–359, 2008.

[10] Carlos Busso et al., "MSP-IMPROV: An acted corpus of dyadic interactions to study emotion perception," *IEEE Transactions on Affective Computing*, vol. 8, no. 1, pp. 67–80, 2016.

[11] Steven R Livingstone and Frank A Russo, "The Ryerson audio-visual database of emotional speech and song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in north american english," *PloS one*, vol. 13, no. 5, pp. e0196391, 2018.

[12] Philip Jackson and SJUoSG Haq, "Surrey audio-visual expressed emotion (savee) database," *University of Surrey: Guildford, UK*, 2014.

[13] Sanyuan Chen et al., "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.

[14] Moataz El Ayadi, Mohamed S Kamel, and Fakhri Karray, "Survey on speech emotion recognition: Features, classification schemes, and databases," *Pattern Recognition*, vol. 44, no. 3, pp. 572–587, 2011.

[15] Tsung-Yi Lin et al., "Focal loss for dense object detection," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.

[16] Yandong Wen et al., "A discriminative feature learning approach for deep face recognition," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*, 2016, pp. 499–515.

[17] Nam Hyeon-Woo et al., "Scratching visual transformer's back with uniform attention," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, pp. 5807–5818.

[18] Prannay Khosla et al., "Supervised contrastive learning," in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020, pp. 18661–18673.

[19] Xiaoyu Tang, Jiazheng Huang, Yixin Lin, Ting Dang, and Jintao Cheng, "Speech emotion recognition via cnn-transformer and multidimensional attention mechanism," *Speech Communication*, p. 103242, 2025.

[20] Chenjing Sun, Yi Zhou, Xin Huang, Jichen Yang, and Xianhua Hou, "Combining wav2vec 2.0 fine-tuning and conlearnnet for speech emotion recognition," *Electronics*, vol. 13, no. 6, pp. 1103, 2024.

[21] Ni Wang and Danyu Yang, "Speech emotion recognition using fine-tuned wav2vec2. 0 and neural controlled differential equations classifier," *PloS one*, vol. 20, no. 2, pp. e0318297, 2025.

[22] Yurun He, Nobuaki Minematsu, and Daisuke Saito, "Multiple acoustic features speech emotion recognition using cross-attention transformer," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.

[23] Yuan Gao, Chenhui Chu, and Tatsuya Kawahara, "Two-stage finetuning of wav2vec 2.0 for speech emotion recognition with ASR and gender pretraining," in *INTERSPEECH*, 2023.

[24] Junyi He et al., "Multilevel transformer for multimodal emotion recognition," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.

[25] Suzhen Wang, Yifeng Ma, and Yu Ding, "Exploring complementary features in multi-modal speech emotion recognition," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.

[26] Lili Guo et al., "Representation learning with spectro-temporal-channel attention for speech emotion recognition," in *ICASSP*, 2021, pp. 6304–6308.

[27] Anish Nediyanchath, Periyasamy Paramasivam, and Promod Yenigalla, "Multi-head attention for speech emotion recognition with auxiliary learning of gender recognition," in *ICASSP*, 2020, pp. 7179–7183.

[28] Mingke Xu et al., "Speech emotion recognition with multiscale area attention and data augmentation," in *ICASSP*, 2021, pp. 6319–6323.

[29] Qi Cao et al., "Hierarchical network based on the fusion of static and dynamic features for speech emotion recognition," in *ICASSP*, 2021, pp. 6334–6338.

[30] Rui Liu et al., "Contrastive learning based modality-invariant feature acquisition for robust multimodal emotion recognition with missing modalities," *IEEE Transactions on Affective Computing*, 2024.

[31] Lu-Yao Liu et al., "ATDA: Attentional temporal dynamic activation for speech emotion recognition," *Knowledge-Based Systems*, vol. 243, pp. 108472, 2022.

[32] Alexei Baevski et al., "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12449–12460, 2020.

[33] Licai Sun et al., "HiCMAE: Hierarchical contrastive masked autoencoder for self-supervised audio-visual emotion recognition," *Information Fusion*, vol. 108, pp. 102382, 2024.

[34] Shaode Yu, Jiajian Meng, Wenqing Fan, Ye Chen, Bing Zhu, Hang Yu, Yaoqin Xie, and Qiurui Sun, "Speech emotion recognition using dual-stream representation and cross-attention fusion," *Electronics*, vol. 13, no. 11, pp. 2191, 2024.

[35] Kateryna Chumachenko, Alexandros Iosifidis, and Moncef Gabbouj, "Self-attention fusion for audiovisual emotion recognition with incomplete data," in *2022 26th International Conference on Pattern Recognition (ICPR)*, 2022, pp. 2822–2828.

[36] Samir Sadok, Simon Leglaive, and Renaud Séguier, "A vector quantized masked autoencoder for speech emotion recognition," in *2023 IEEE International Conference on Acoustics, Speech, and Signal processing Workshops (ICASSPW)*, 2023, pp. 1–5.

[37] Li-Wei Chen and Alexander Rudnicky, "Exploring wav2vec 2.0 fine tuning for improved speech emotion recognition," in *ICASSP*, 2023, pp. 1–5.