



A Multi-Level Alignment and Cross Modal Unified Semantic Graph Refinement Network for Conversational Emotion Recognition

Xiaoheng Zhang, Weigang Cui, Bin Hu, Fellow, IEEE, and Yang Li, Senior, IEEE

Code:<https://github.com/zxiaohen/MA-CMU-SGRNet>

— IEEE Transactions on Affective Computing 2024

2024. 3. 17 • ChongQing



gesis
Leibniz-Institut
für Sozialwissenschaften



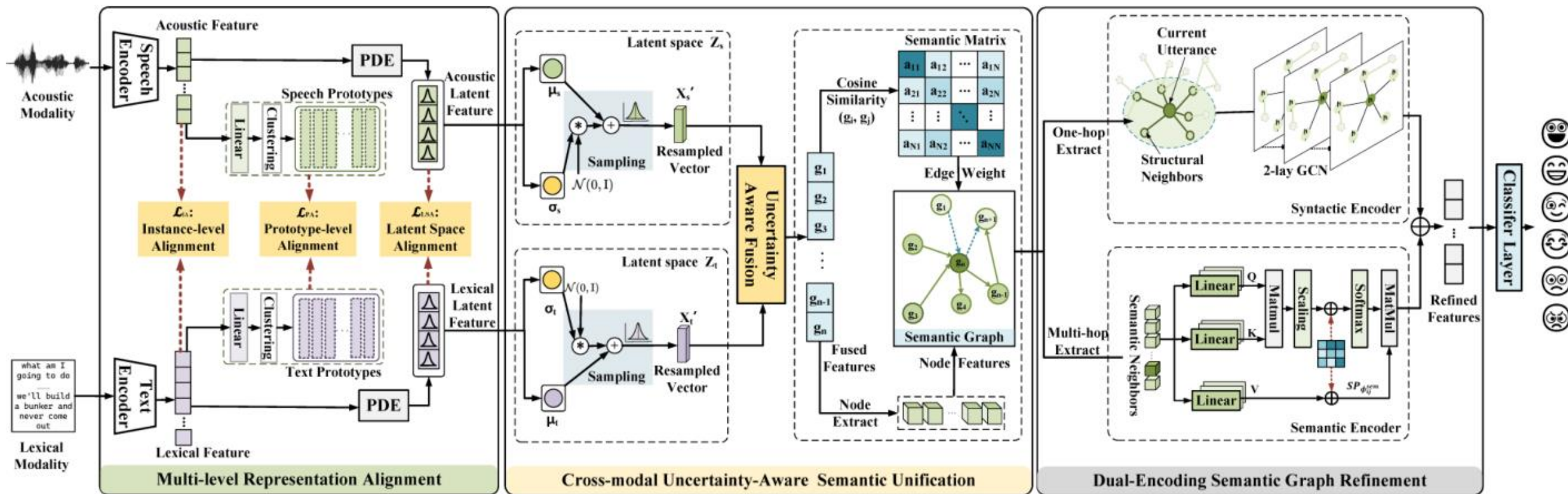
Reported by JiaWei Cheng



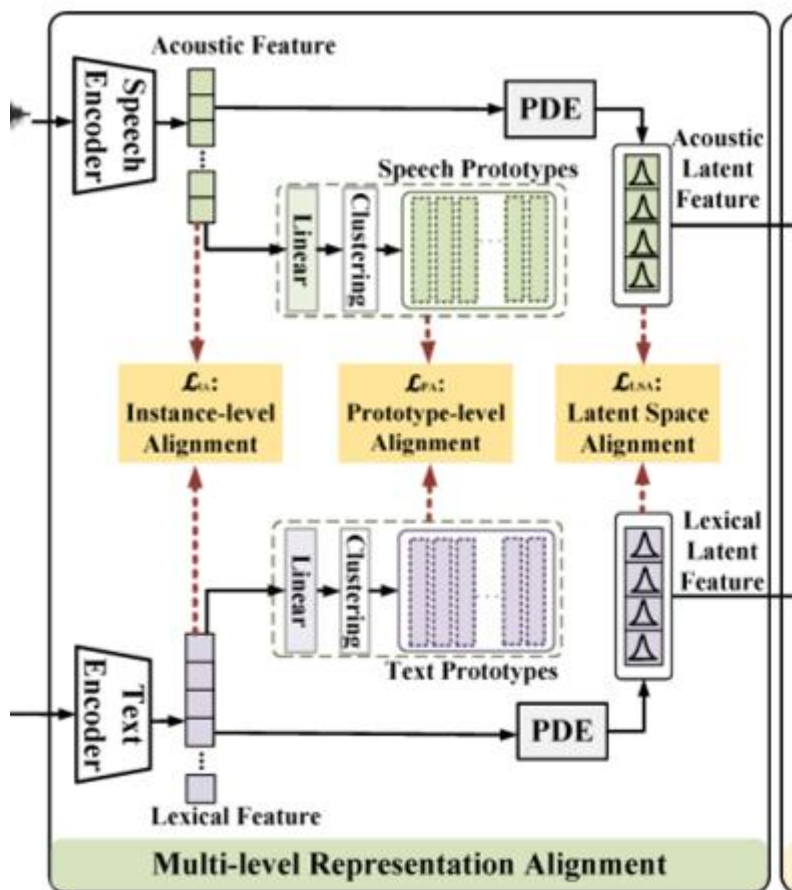
Motivation

- (1) The method of supervised contrastive learning is too tough.
- (2) Previous fusion methods did not consider the uncertainty in each modal
- (3) Previous methods did not fully mine the semantic context information in the conversation

Overview



Method



$$S_{i,j}^{(0)} = \exp\left(\frac{1}{\tau_1} \cdot \frac{z_i^T z_j}{\|z_i\| \|z_j\|}\right), z \in \{\tilde{s}, \tilde{t}\} \quad (1)$$

$$S_{i,j}^{(1)} = \exp\left(\frac{1}{\tau_1} \cdot \frac{\tilde{z}_i^T z_j}{\|\tilde{z}_i\| \|z_j\|} - \xi\right), z \in \{\tilde{s}, \tilde{t}\} \quad (2)$$

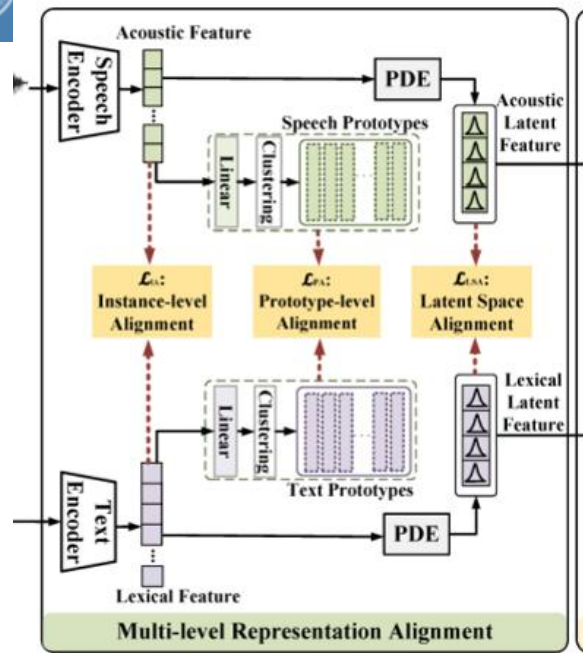
$$S_{i,j} = \exp\left(\frac{1}{\tau_{m(i,j)}} \cdot \frac{z_i^T z_j}{\|z_i\| \|z_j\|} - \xi_{m(i,j)}\right), z \in \{\tilde{s}, \tilde{t}\} \quad (3)$$

$$l_i^{s2t} = -\frac{1}{|P(i)|} \sum_{j \in P(i)} \log \frac{S_{i,j}}{S_{i,j} + \sum_{n \in N(i)} S_{i,n}} \quad (4)$$

$$l_i^{t2s} = -\frac{1}{|P(i)|} \sum_{j \in P(i)} \log \frac{S_{j,i}}{S_{j,i} + \sum_{n \in N(i)} S_{n,i}} \quad (5)$$

$$j \in P(i) = \{j | j \in \mathcal{B}, y_j = y_i, S_{i,j} > 1\}$$

$$\mathcal{L}_{IA} = \frac{1}{2N_p} \sum_{i=1}^{N_p} (l_i^{s2t} + l_i^{t2s}) \quad (6)$$



Method

$$p_{s,i}^k = \frac{\exp(\tilde{s}_i^T c_k / \tau_2)}{\sum_k \exp(\tilde{s}_i^T c_k / \tau_2)} \quad (7)$$

$$p_{t,i}^k = \frac{\exp(\tilde{t}_i^T c_k / \tau_2)}{\sum_k \exp(\tilde{t}_i^T c_k / \tau_2)} \quad (8)$$

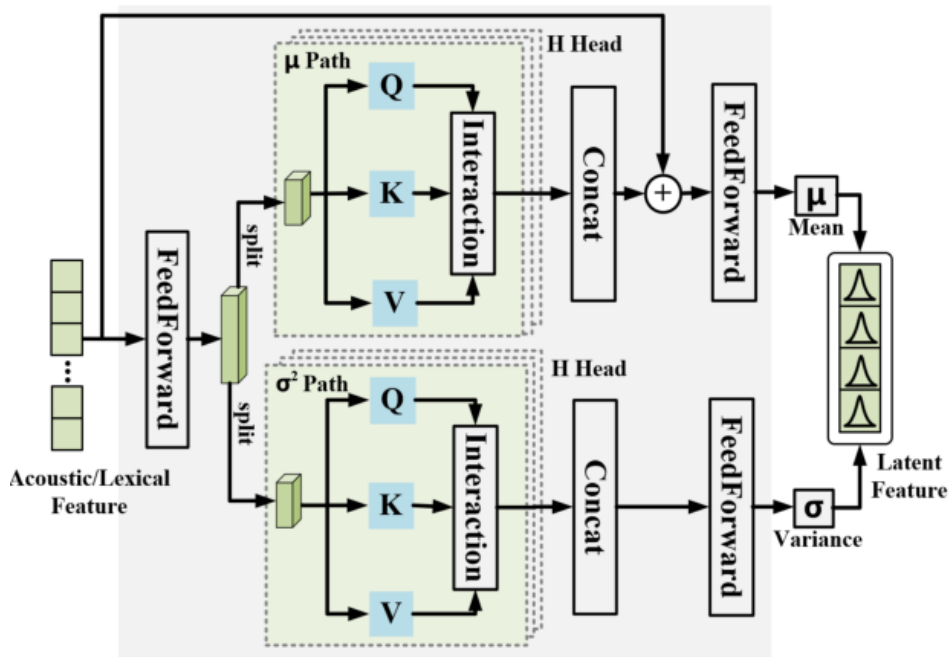
$$l(\tilde{s}_i, q_{t,j}) = \sum_{k=1}^K q_{t,i}^{(k)} \log p_{s,i}^{(k)} \quad (9)$$

$$l(\tilde{t}_i, q_{s,j}) = \sum_{k=1}^K q_{s,i}^{(k)} \log p_{t,i}^{(k)} \quad (10)$$

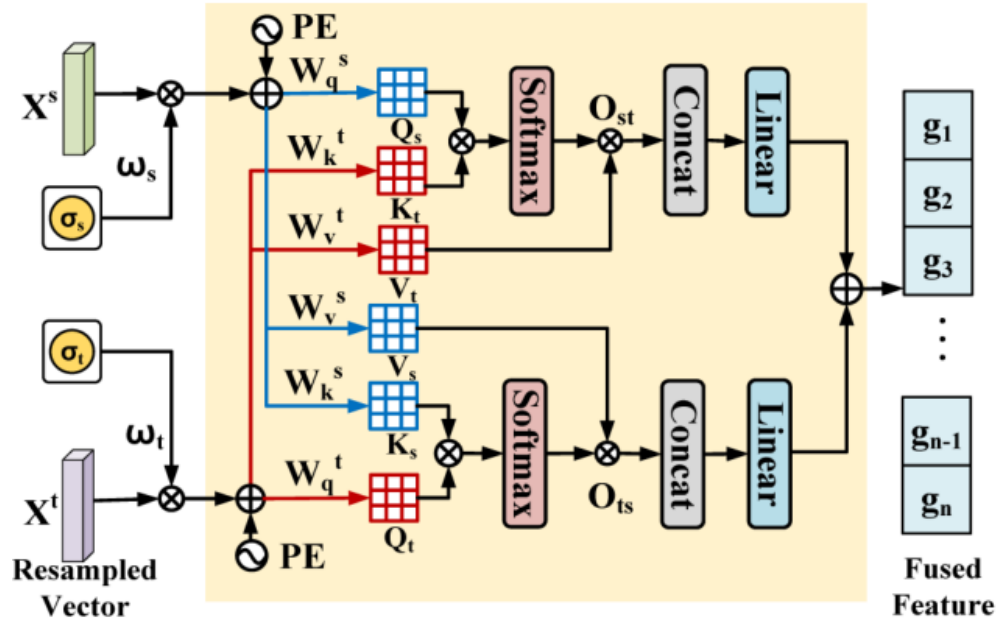
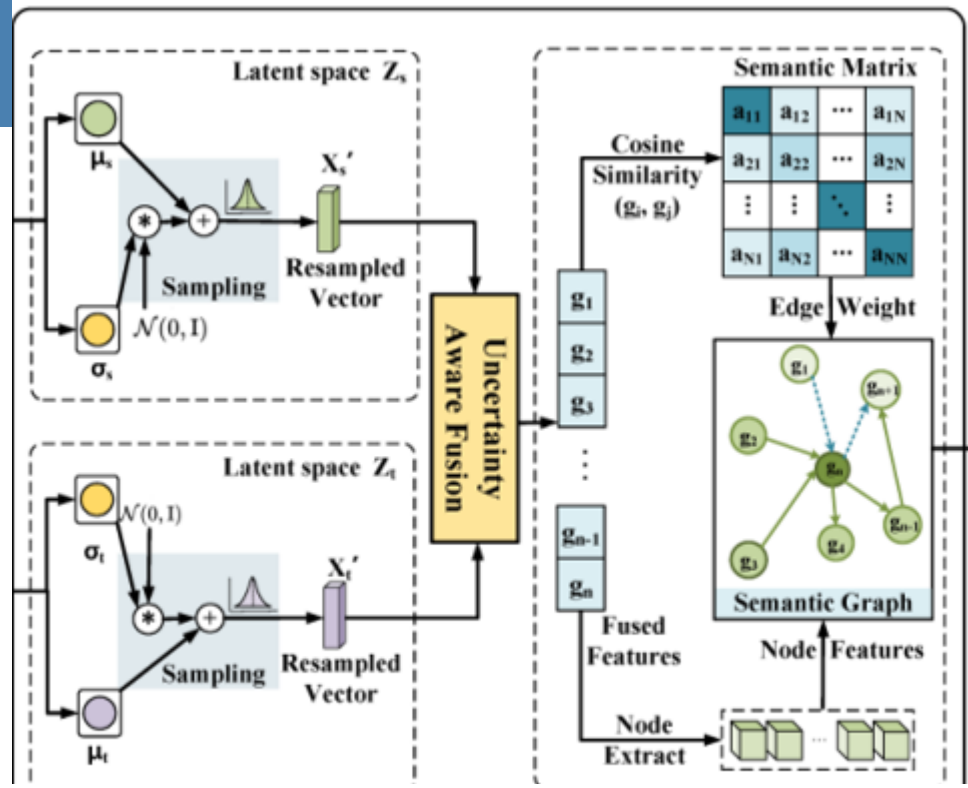
$$\mathcal{L}_{PA} = \frac{1}{2N_p} \sum_{i=1}^{N_p} (l(\tilde{s}_i, q_{t,i}) + l(\tilde{t}_i, q_{s,i})) \quad (11)$$

$$\mathcal{L}_{LSA} = \frac{1}{2N_p} \sum_{i=1}^{N_p} (\|\mu_{s,i} - \mu_{t,i}\|_2^2 + \|\sigma_{s,i} - \sigma_{t,i}\|_2^2)^{\frac{1}{2}} \quad (12)$$

$$\mathcal{L}_A = \gamma_1 \mathcal{L}_{IA} + \gamma_2 \mathcal{L}_{PA} + \gamma_3 \mathcal{L}_{LSA} \quad (13)$$



Method



$$\omega_s = \frac{\|\sigma_s\|}{\|\sigma_s\| + \|\sigma_t\|} \quad (14)$$

$$\omega_t = \frac{\|\sigma_t\|}{\|\sigma_s\| + \|\sigma_t\|} \quad (15)$$

$$\tilde{X}^s = \omega_s \odot \check{X}^s \quad (16)$$

$$\tilde{X}^t = \omega_t \odot X^t \quad (17)$$

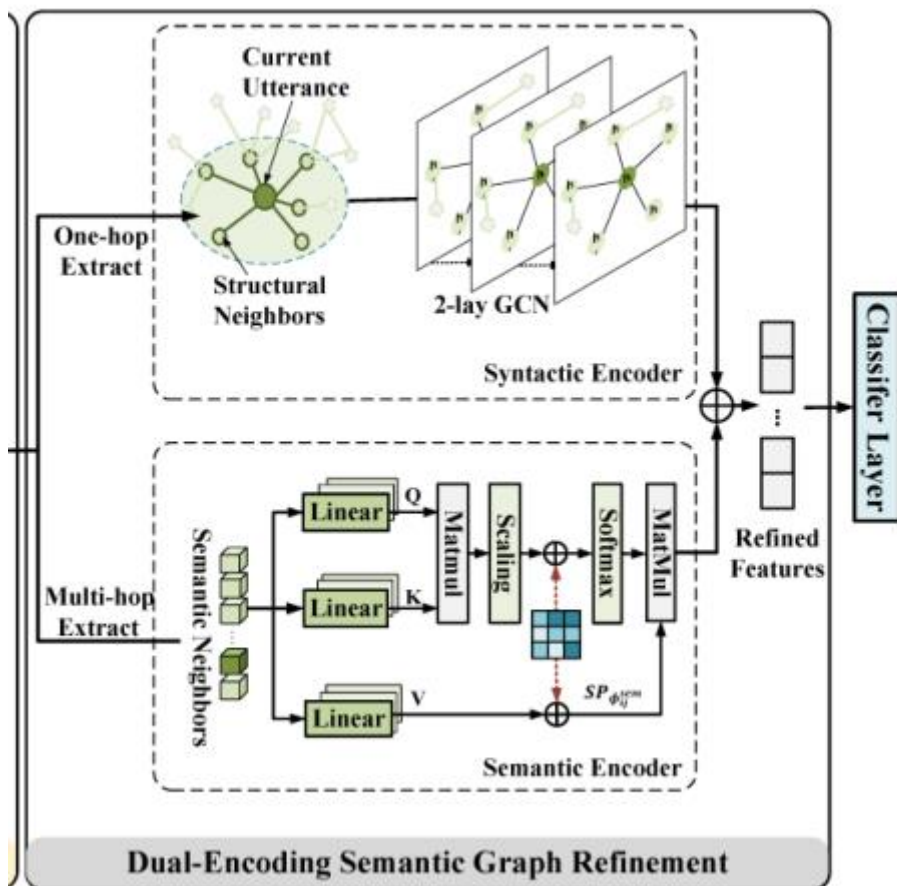
$$Q_s = \tilde{X}^s W_q^s \quad (18)$$

$$K_t = \tilde{X}^t W_k^t \quad (19)$$

$$V_t = \tilde{X}^t W_v^t \quad (20)$$

$$O_{st} = \text{softmax}(Q_s^T K_t) \cdot V_t \quad (21)$$

Method



$$\mathcal{M}_{i,j}^s = 1 - \arccos\left(\frac{g_i^T g_j}{\|g_i\| \|g_j\|}\right), \quad i, j \in [1, N] \quad (22)$$

$$H_i^{(l+1)} = F\left(\sum_{r \in R} \sum_{j \in G_i^r} \frac{1}{|G_i^r|} W_r^{(l)} g_j^{(l)} + W_0^{(l)} g_i^{(l)}\right) \quad (23)$$

$$SP_{\phi_{ij}^{sem}} = \mathcal{M}_{i,j}^s + \mathcal{M}_{i,j}^p \quad (24)$$

$$a_{i,j} = \frac{(W_q g_i)^T (W_k g_j)}{\sqrt{d_f}} + SP_{\phi_{ij}^{sem}} \quad (25)$$

$$h_i^{se} = \sum_{j=1}^{N_0} \text{softmax}(a_{i,j}) (W_v g_j + SP_{\phi_{ij}^{sem}}) \quad (26)$$

$$\tilde{h}_i = \text{ReLU}(W_f h_i + b_f) \quad (27)$$

$$\beta_i = \text{softmax}(W_p \tilde{h}_i + b_p) \quad (28)$$

$$\hat{y}_i = \text{argmax}(\beta_i) \quad (29)$$

$$\mathcal{L}_{cl} = -\frac{1}{\sum_{i=1}^L N_i} \sum_{i=1}^L \sum_{c=1}^{C_0} y_{i,c}^{(j)} \cdot \log \hat{y}_{i,c}^{(j)} \quad (30)$$

$$\mathcal{L}_{total} = \mathcal{L}_{cl} + \lambda \mathcal{L}_A \quad (31)$$

Experiments

TABLE I: COMPARISON WITH THE STATE-OF-THE-ART METHODS ON IEMOCAP DATASET.

Models	Year	IEMOCAP: Emotion Categories													
		Happy		Sad		Neutral		Angry		Excited		Frustrated		Average	
		Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	WAA	WAF1
Bc-LSTM [44]	2017	22.5	35.6	58.6	69.2	56.5	53.5	70.0	66.3	58.8	61.1	67.4	62.4	59.8	59.0
DialogueRNN [12]	2019	31.25	33.8	66.1	69.8	63.0	57.7	61.7	62.5	61.5	64.4	59.6	59.5	59.3	59.8
CTNet [3]	2021	47.9	51.3	78.0	79.9	69.0	65.8	72.9	67.2	85.3	78.7	52.2	58.8	68.0	67.5
A-DMN [11]	2022	43.1	50.6	69.4	76.8	63.0	62.9	63.5	56.5	88.3	77.9	53.3	55.7	64.6	64.3
I-GCN [14]	2022	51.4	50.0	85.3	83.8	60.4	59.3	61.2	64.6	75.6	74.3	57.2	59.0	65.5	65.4
GraphCFC [16]	2023	-	43.1	-	85.0	-	64.7	-	71.4	-	78.9	-	63.7	-	68.9
Ours	2023	52.6	57.1	78.8	79.9	74.3	71.0	75.2	71.5	80.3	78.4	65.1	67.5	72.4	71.6

The improvement is statistically significant with $p \leq 0.05$ under t -test. Bold font represents the best performance. Acc. = Accuracy.

Experiments

TABLE II: COMPARISON WITH THE STATE-OF-THE-ART METHODS ON MELD DATASET.

Models	Year	MELD: Emotion Categories							
		Anger	Disgust	Fear	Joy	Neutral	Sadness	Surprise	Avg
		F1	F1	F1	F1	F1	F1	F1	WAF1
Bc-LSTM [44]	2017	43.4	23.7	9.4	54.5	76.7	24.3	51.0	59.3
DialogueRNN [12]	2019	43.7	7.9	11.7	54.4	77.4	34.6	52.5	60.3
CTNet [3]	2021	44.6	11.2	10.0	56.0	77.4	32.5	52.7	60.5
A-DMN [11]	2022	43.9	7.2	12.0	56.7	77.1	29.1	55.1	60.4
I-GCN [14]	2022	43.5	11.8	8.0	54.7	78.0	38.5	51.6	60.8
CMCF-SRNet [15]	2023	43.9	10.9	11.5	55.8	77.2	36.0	52.9	61.8
Ours	2023	44.3	11.9	12.1	56.9	78.4	35.9	53.5	62.3



Experiments

TABLE III PERFORMANCE ON CMU-MOSEI DATASET.

Methods	CMU-MOSEI		
	Year	WAA	WAF1
GMFN [42]	2017	76.9	77.0
MuT [48]	2019	82.5	82.3
MMIM [49]	2021	85.9	85.9
CONKI [50]	2023	86.2	86.1
Ours	2023	86.9	86.8

Experiments

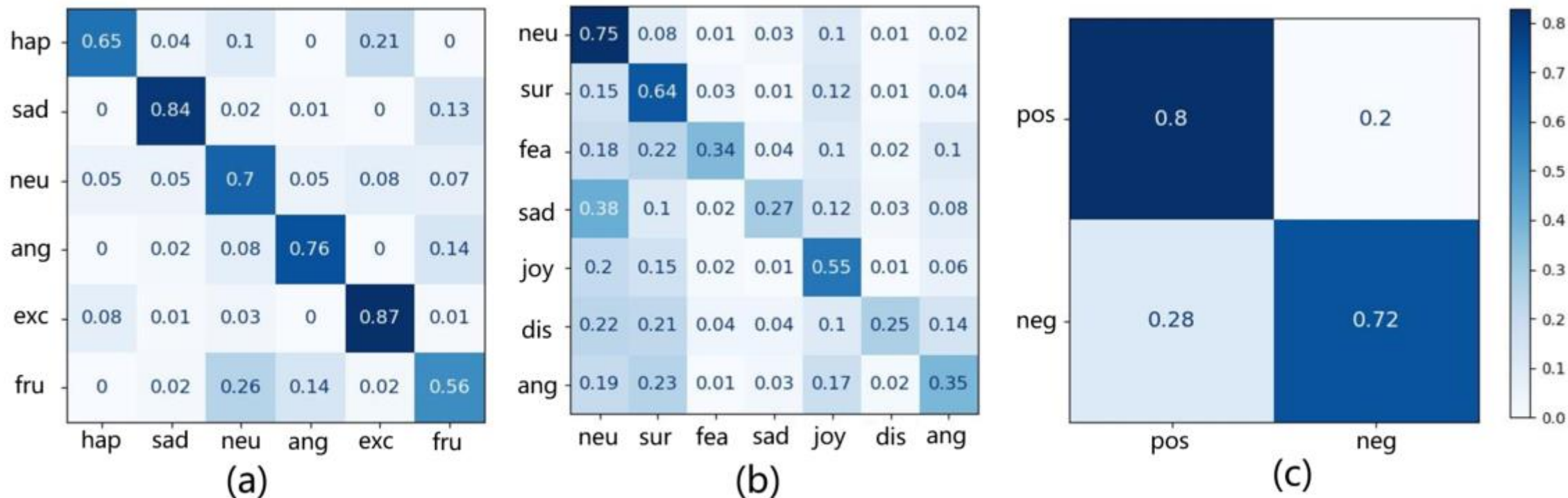


Fig. 4. The confusion matrices of the proposed MA-CMU-SGRNet on three datasets. (a) IEMOCAP (b) MELD (c) CMU-MOSEI

Experiments

TABLE IV
RESULTS OF ABLATION STUDIES ON THREE DATASETS.

Methods	IEMOCAP		MELD		CMU-MOSEI	
	WAA	WAF1	WAA	WAF1	WAA	WAF1
w/o UW	68.9 \pm 0.48 [†]	67.6 \pm 0.61 [†]	60.3 \pm 0.34 [†]	59.7 \pm 0.75 [†]	83.7 \pm 0.73 [†]	84.1 \pm 0.52 [*]
w/o CMI	68.1 \pm 0.61 [†]	67.4 \pm 0.35 [†]	60.5 \pm 0.73 [†]	59.9 \pm 0.42 [*]	82.8 \pm 0.54 [†]	83.2 \pm 0.42 [†]
w/o SyE	68.3 \pm 0.73 [†]	66.4 \pm 0.54 [†]	61.7 \pm 0.64 [*]	61.4 \pm 0.52 [*]	83.2 \pm 0.45 [†]	83.6 \pm 0.56 [*]
w/o SeE	68.8 \pm 0.53 [†]	67.9 \pm 0.67 [†]	59.5 \pm 0.32 [†]	59.2 \pm 0.47 [†]	82.9 \pm 0.54 [†]	83.1 \pm 0.65 [†]
Ours	72.4\pm0.61	71.6\pm0.73	62.8\pm0.54	62.3\pm0.62	86.9\pm0.66	86.8\pm0.92

where the symbols [†] and * indicate that the difference with respect to the ablation setting is statistically significant at $p < 0.001$ [†] and $p < 0.01$ ^{*} under t -test.

Experiments

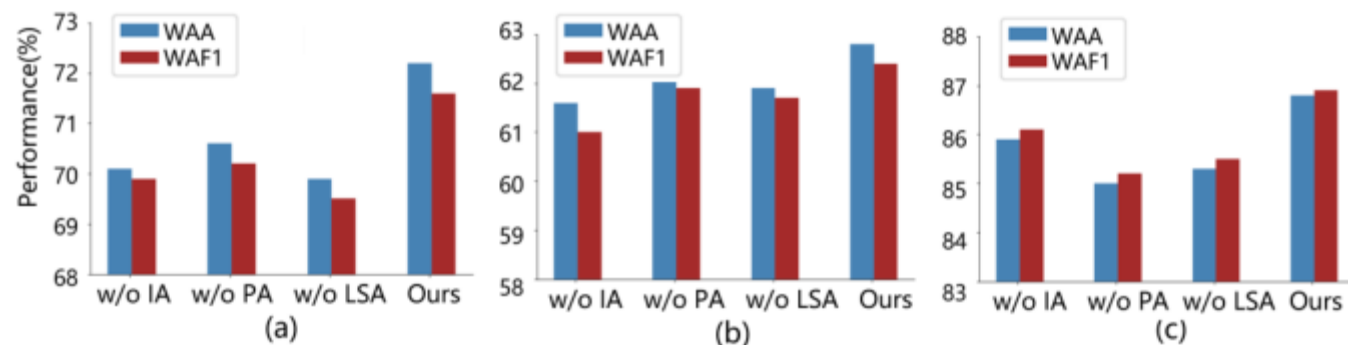


Fig. 5. Comparison of WAA and WAF1 on three datasets.
(a) IEMOCAP (b) MELD (c) CMU-MOSEI.

Experiments

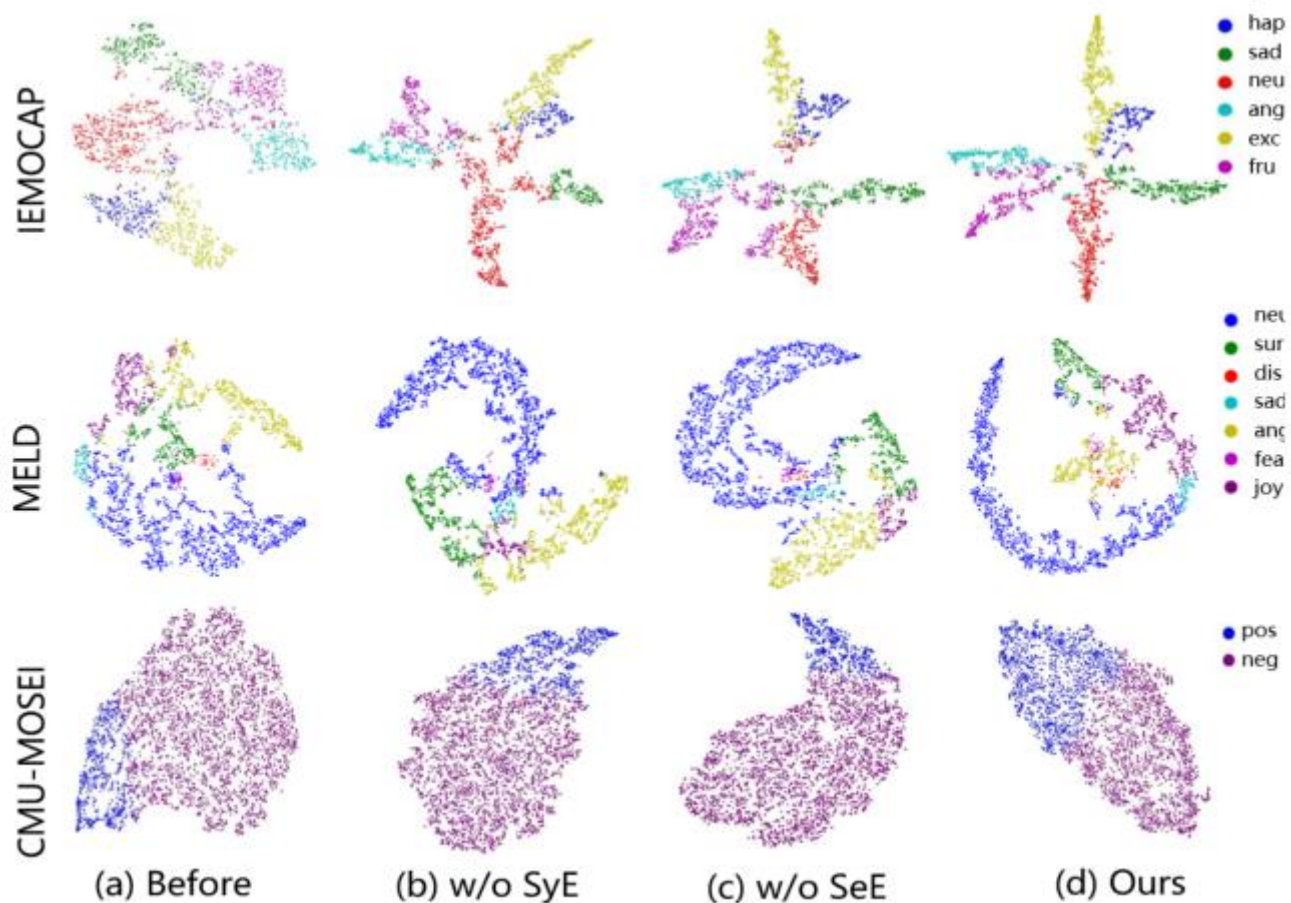


Fig. 7. The visualization of the t-SNE representations

Experiments

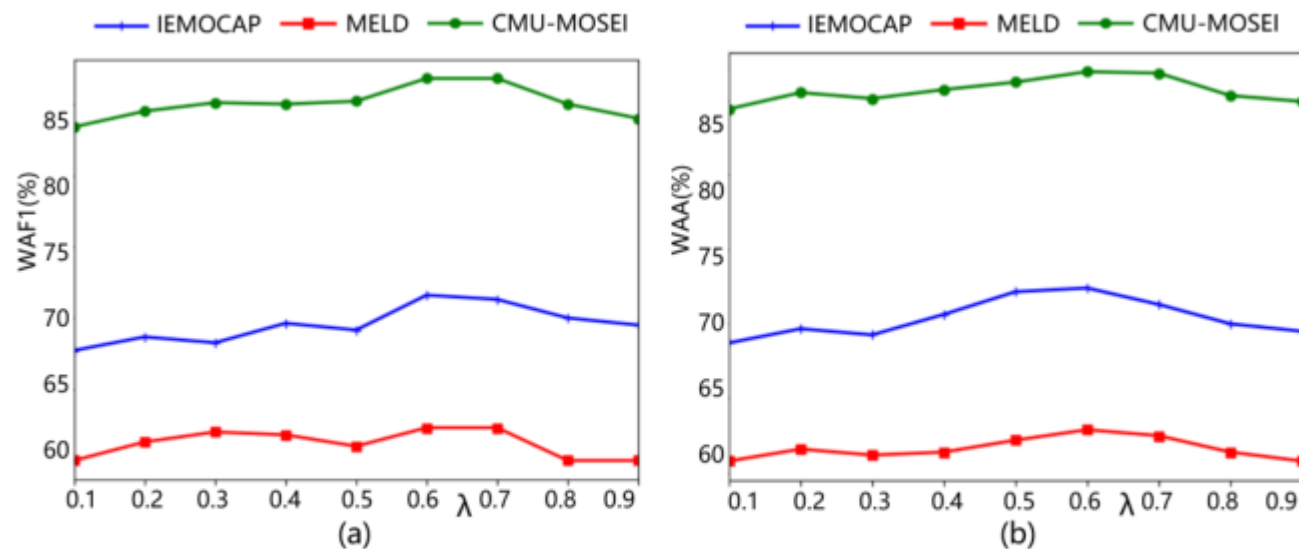


Fig. 8. Comparison of WAA and WAF1 on three datasets.
(a) IEMOCAP (b) MELD (c) CMU-MOSEI.



Thanks!