



GraphCFC: A Directed Graph based Cross-modal Feature Complementation Approach for Multimodal Conversational Emotion Recognition

Jiang Li, Xiaoping Wang, *Senior Member, IEEE*, Guoqing Lv, and Zhigang Zeng, *Fellow, IEEE*

Code:None

— IEEE Transactions on Multimedia 2023

2024. 2. 1 • ChongQing

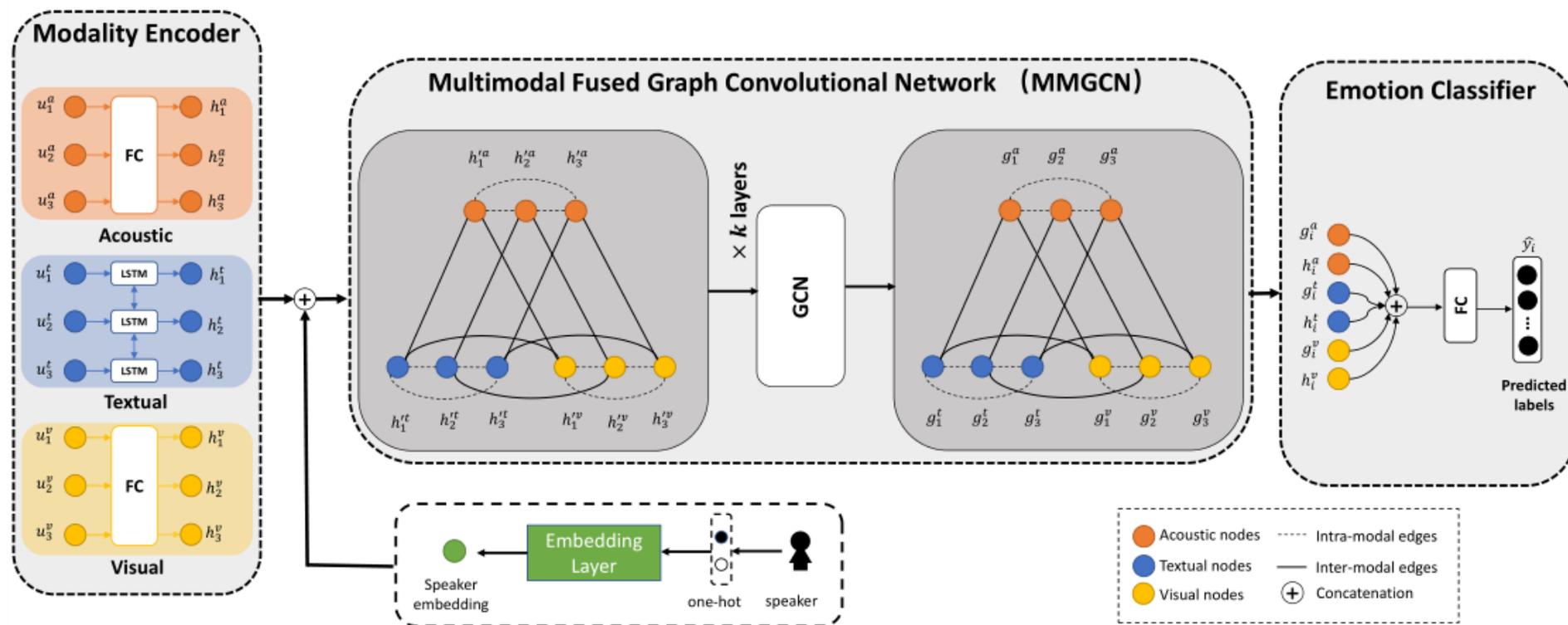


gesis
Leibniz-Institut
für Sozialwissenschaften



Reported by JiaWei Cheng

Motivation



MMGCN: Multimodal Fusion via Deep Graph Convolution Network for Emotion Recognition in Conversation(ACL2021)

This approach not only adds redundant information due to inconsistent data distribution among modalities, but also may risk losing diverse information in the conversational graph

Overview

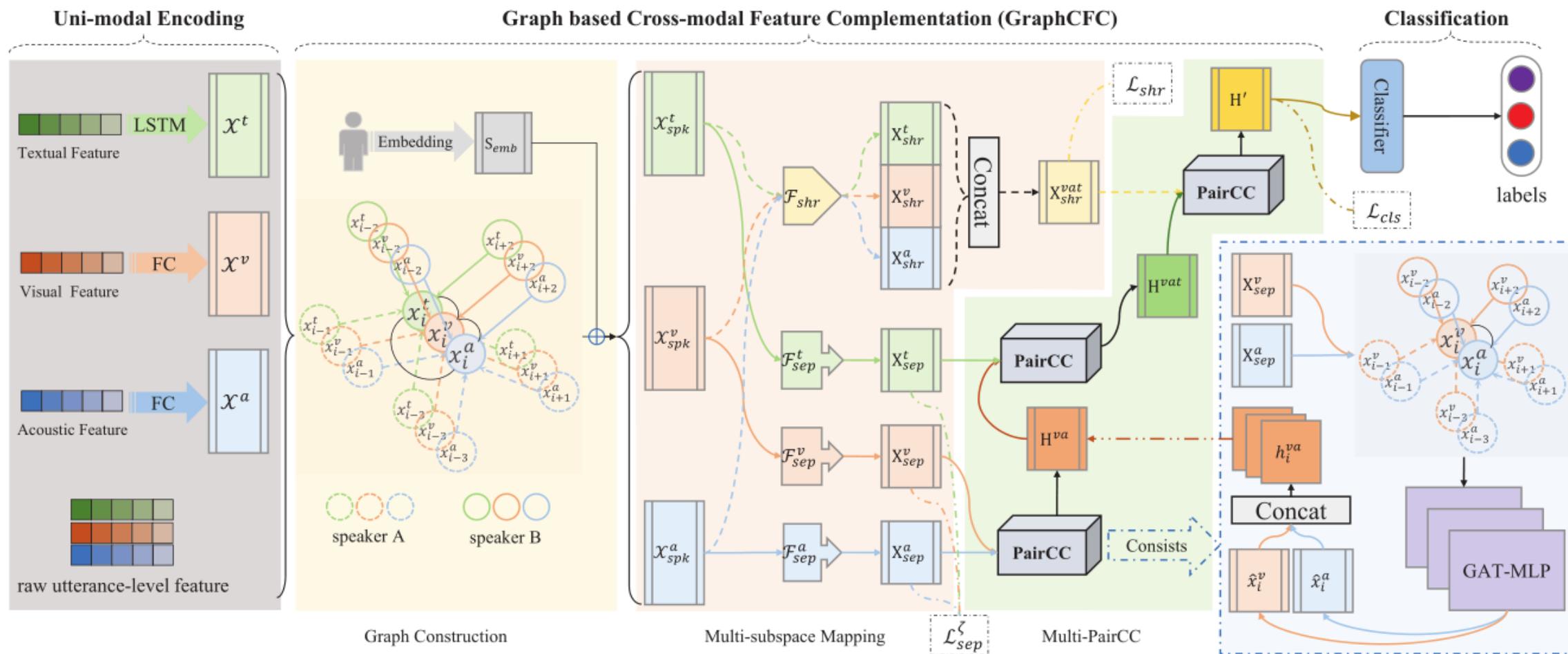
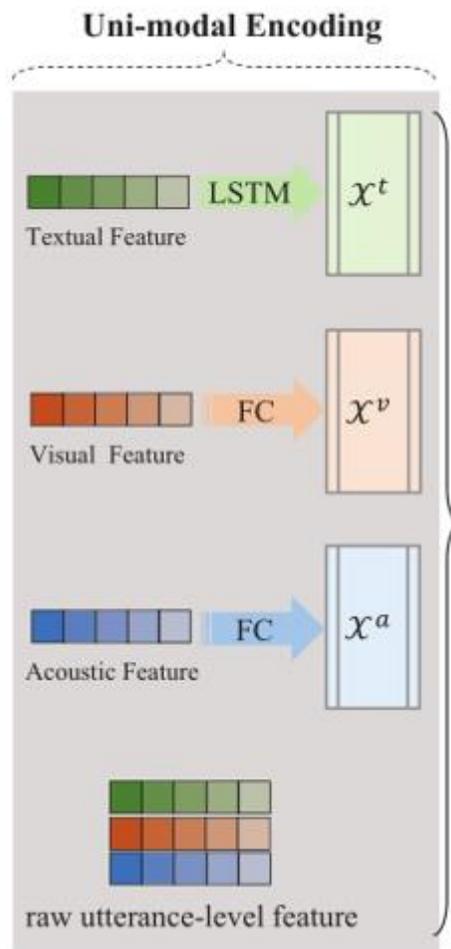


Fig. 2: The illustration of graph-based multimodal ERC, which includes uni-modal encoding, graph based cross-modal feature completion and multimodal emotion classification.

Method



$$x_i^t, x_{h,i}^t = \overrightarrow{\text{LSTM}}(u_i^t; \Theta_{ls}^t), \quad (1)$$

$$x_i^\tau = \text{FC}(u_i^\tau; \Theta_{fc}^\tau), \tau \in \{a, v\}, \quad (2)$$

Method

$$\mathcal{E}_{intra} = \begin{cases} \{(u_t^P, u_i^P) | i - j < t < i - 1\} \\ \{(u_i^P, u_t^P) | i + 1 < t < i + k\} \end{cases}, \quad (3)$$

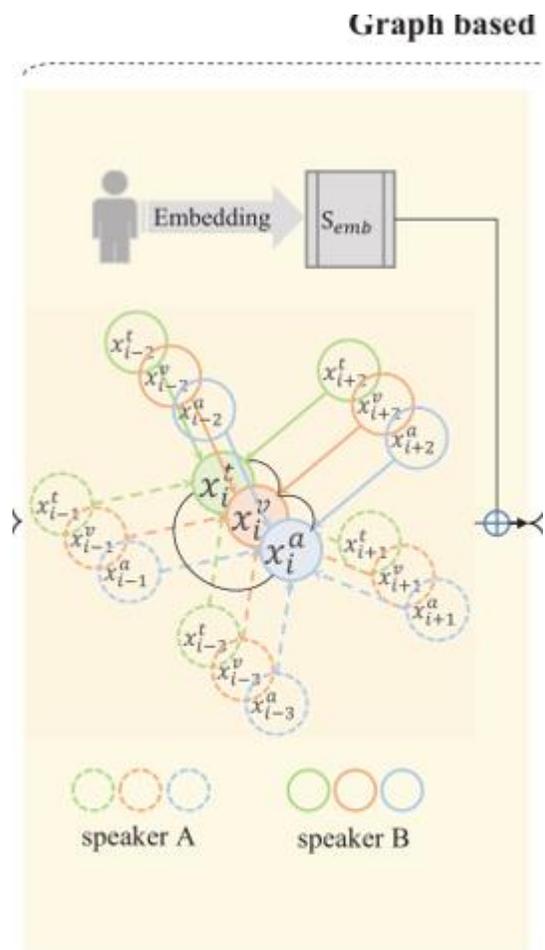
$$\mathcal{E}_{inter} = \{(u_i^P, u_i^Q), (u_i^Q, u_i^P)\}, \quad (4)$$

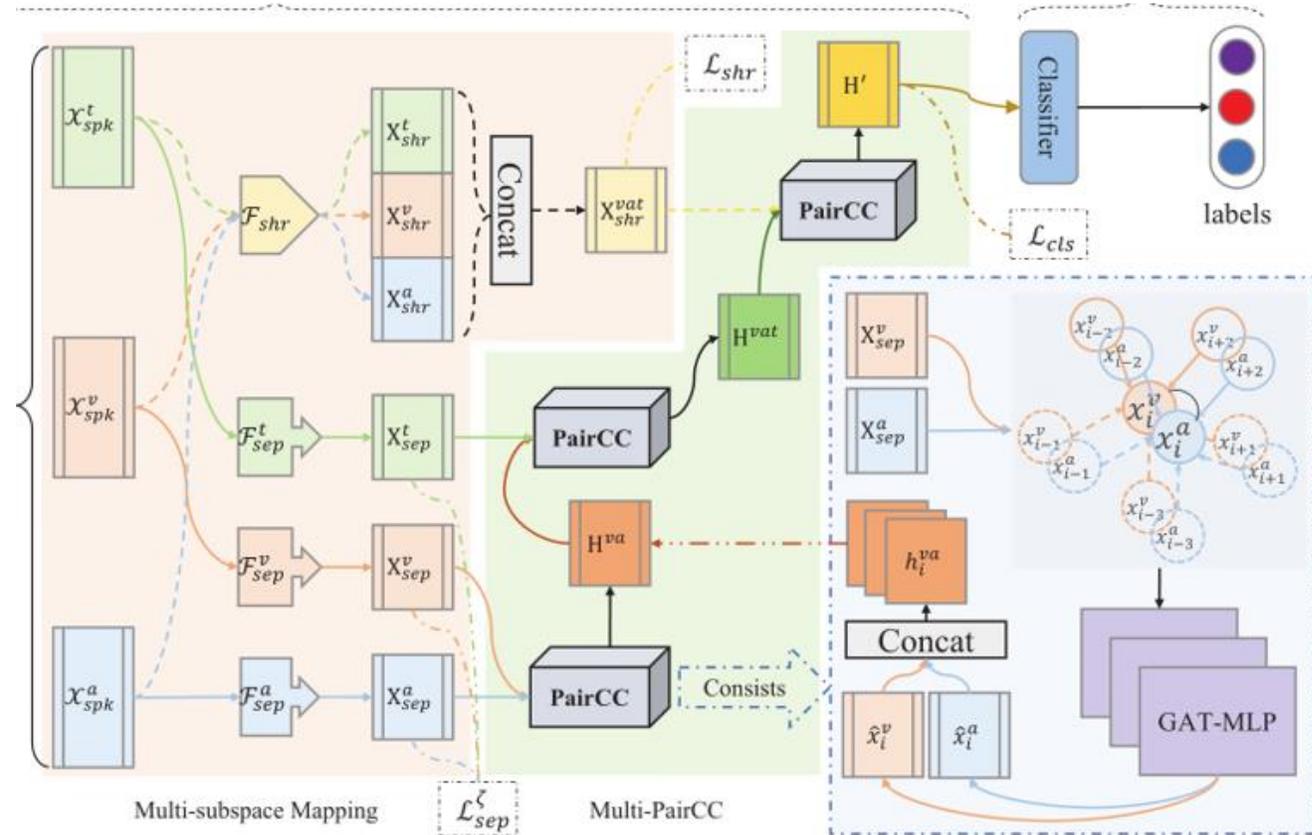
$$ET_{intra} = \{et(s_1, s_1), et(s_1, s_2), et(s_1, s_3), et(s_2, s_2), et(s_2, s_3), et(s_3, s_3)\}. \quad (5)$$

$$ET_{inter} = \{et(mod_1, mod_2), et(mod_1, mod_3), et(mod_2, mod_3)\}. \quad (6)$$

$$S_{emb} = \text{Embedding}(S, D), \quad (7)$$

$$\mathcal{X}_{spk}^\zeta = \mu S_{emb} + \mathcal{X}^\zeta, \quad (8)$$





$$\begin{aligned} X_{shr}^z &= \mathcal{F}_{shr}(X_{spk}^z; \Theta_{shr}), \\ X_{shr}^{vat} &= \text{Lin}([X_{shr}^v \parallel X_{shr}^a \parallel X_{shr}^t]; \Theta'_{shr}), \end{aligned} \quad (9)$$

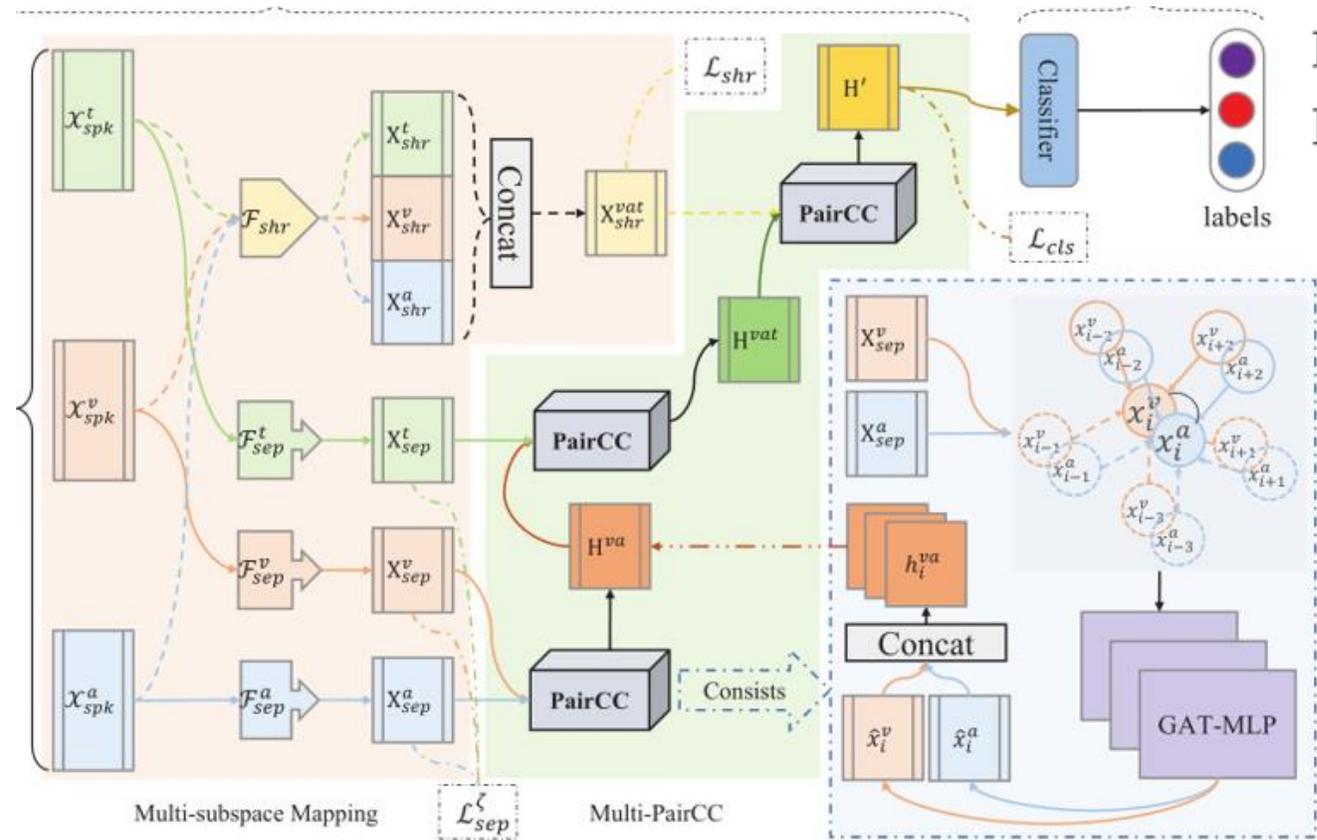
$$\begin{aligned} X_{sep}^z &= \mathcal{F}_{sep}(X_{spk}^z; \Theta_{sep}^z), \\ \mathcal{F}(X; \Theta) &= \text{Norm}(\text{Drop}(\text{Lin}(\text{Drop}(\sigma(\text{Lin}(X; \Theta_0))))); \Theta_1)), \end{aligned} \quad (10)$$

$$\begin{aligned} v'_i &= \text{ReLU}(W'_0 x_{shr,i}^{vat} + b'_0), \\ p'_i &= \text{Softmax}(W'_1 v'_i + b'_1), \end{aligned} \quad (11)$$

$$\mathcal{L}_{shr} = -\frac{1}{\sum_{k=0}^{N-1} n(k)} \sum_{i=0}^{N-1} \sum_{j=0}^{n(i)-1} y_{ij} \log p'_{ij} + \lambda |\Theta'_{re}|, \quad (12)$$

$$\begin{aligned} v_i^z &= \text{ReLU}(W_0^z x_{sep,i}^z + b_0^z), \\ p_i^z &= \text{Softmax}(W_1^z v_i^z + b_1^z), \end{aligned} \quad (13)$$

$$\mathcal{L}_{sep}^z = -\frac{1}{\sum_{k=0}^{N-1} n(k)} \sum_{i=0}^{N-1} \sum_{j=0}^{n(i)-1} y_{ij} \log p_{ij}^z + \lambda |\Theta_{re}^z|, \quad (14)$$



$$\begin{aligned} X_{\text{gat}} &= \text{Norm}(\text{MultiGAT}(\mathcal{E}, X_{\text{in}}; \Theta_{\text{gat}}) + X_{\text{in}}), \\ X_{\text{out}} &= \text{Norm}(\text{FeedForward}(X_{\text{gat}}; \Theta_{\text{fed}}) + X_{\text{gat}}), \end{aligned} \quad (15)$$

$$\begin{aligned} \text{FeedForward}(X_{\text{gat}}, \Theta_{\text{fed}}) &= \\ \text{Drop}(\text{Lin}(\text{Drop}(\sigma(\text{Lin}(X_{\text{gat}}; \Theta_0))); \Theta_1)), \end{aligned} \quad (16)$$

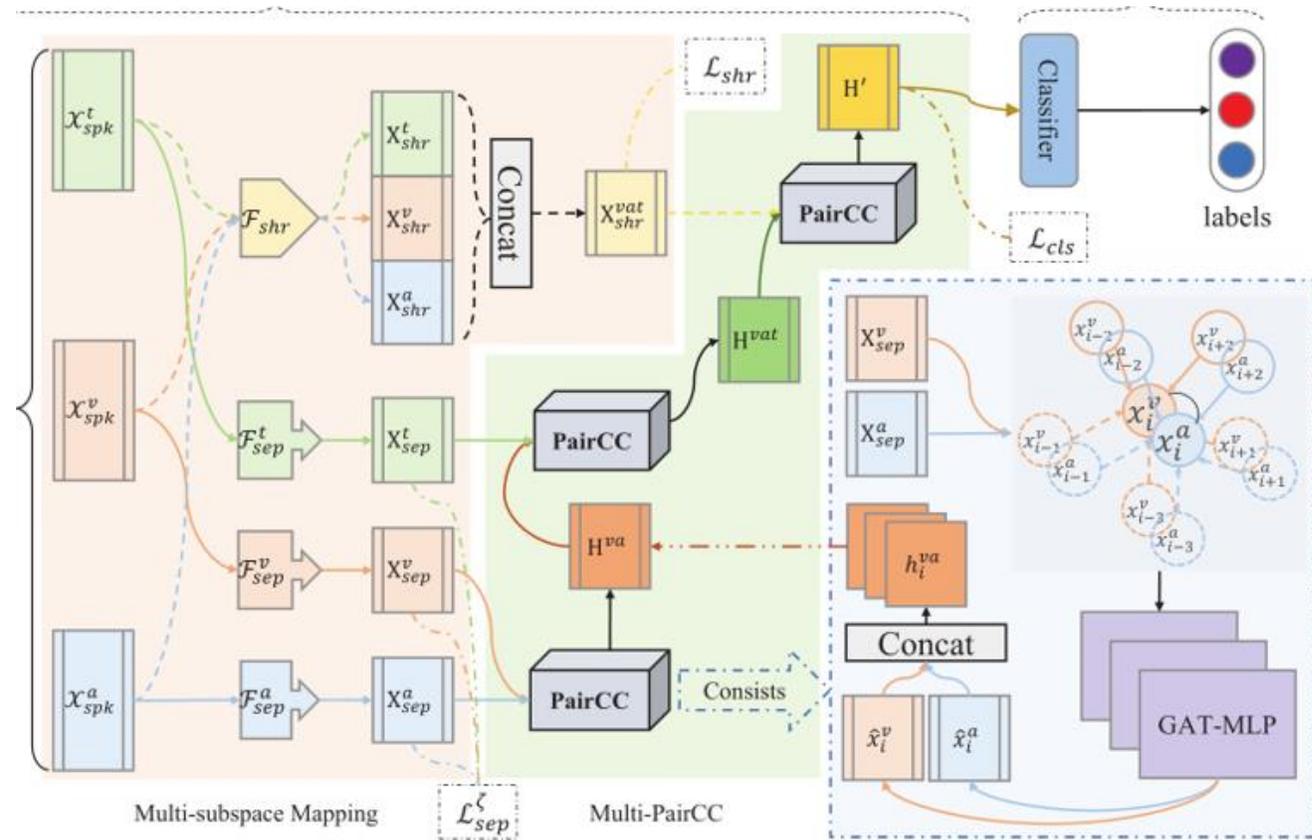
$$\begin{aligned} \text{MultiGAT}(\mathcal{E}, X_{\text{in}}; \Theta_h) &= \Theta_h[\text{head}_1 \parallel \dots \parallel \text{head}_h], \\ \text{where } \text{head}_i &= \text{SingleGAT}(\mathcal{E}, X_{\text{in}}; \Theta_i), \end{aligned} \quad (17)$$

$$\begin{aligned} X_{\text{gat}} &= \text{MultiGAT}(\mathcal{E}, \text{Norm}(X_{\text{in}}); \Theta_{\text{gat}}) + X_{\text{in}}, \\ X_{\text{out}} &= \text{FeedForward}(\text{Norm}(X_{\text{gat}}); \Theta_{\text{fed}}) + X_{\text{gat}}. \end{aligned} \quad (18)$$

$$\begin{aligned} H^{va} &= \text{PairCC}(X_{\text{sep}}^v, X_{\text{sep}}^a; \Theta_{\text{sep}}^{va}), \\ H^{vat} &= \text{PairCC}(X_{\text{sep}}^t, H^{va}; \Theta_{\text{sep}}^{vat}), \end{aligned} \quad (19)$$

$$\begin{aligned} H' &= \text{PairCC}(X_{\text{shr}}^{vat}, H^{vat}; \Theta'), \\ x_{\text{agg},i} &= \text{AGG}(\{x_j | w_j \in \mathcal{N}(w_i)\}; \Theta_{\text{agg}}), \end{aligned} \quad (20)$$

$$\begin{aligned} x_{\text{com},i} &= \text{COM}(x_i, x_{\text{agg},i}; \Theta_{\text{com}}), \\ x_{\text{agg},i} &= \sum_{w_j \in \mathcal{N}(w_i)} \alpha_{ij} W_{\text{agg}} x_j, \end{aligned} \quad (21)$$



$$\alpha_{ij} = \frac{\exp(a^\top \sigma(\Theta_{\text{att}}[x_i \| x_j]))}{\sum_{w_k \in \mathcal{N}(w_i)} \exp(a^\top \sigma(\Theta_{\text{att}}[x_i \| x_k]))}, \quad (22)$$

$$ET_{\text{emb}} = \text{Embedding}(ET, DM), \quad (23)$$

where $DM = M \times (D^2 + D + M - 1)/2$,

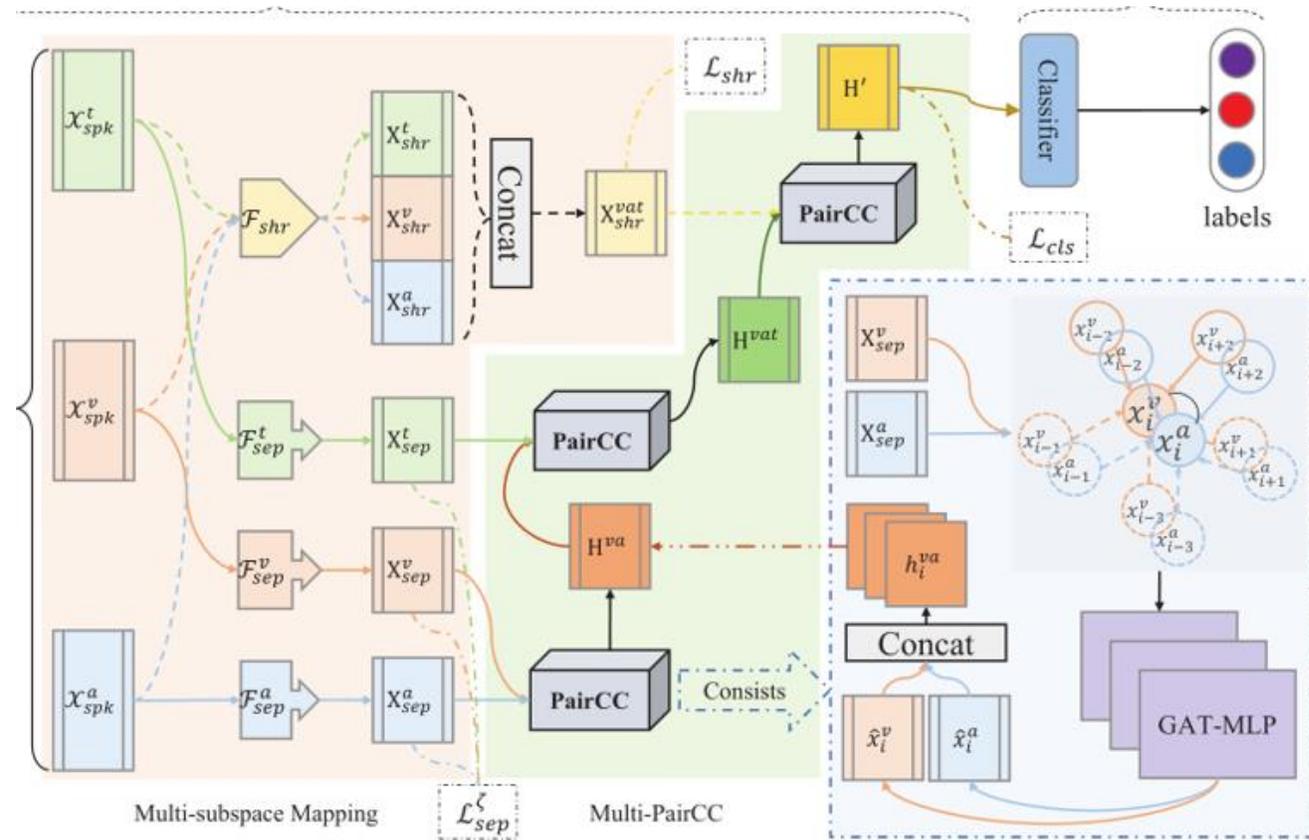
$$\alpha'_{ij} = \frac{\exp(a^\top \sigma(\Theta_{\text{att}}[x_i \| x_j \| et_{ij}]))}{\sum_{w_k \in \mathcal{N}(w_i)} \exp(a^\top \sigma(\Theta_{\text{att}}[x_i \| x_k \| et_{ik}]))}, \quad (24)$$

$$x_{\text{com},i}^1 = \text{GRU}(x_i, x_{\text{agg},i}; \Theta_{\text{com}}^1), \quad (25)$$

$$x_{\text{com},i}^2 = \text{GRU}(x_{\text{agg},i}, x_i; \Theta_{\text{com}}^2), \quad (26)$$

$$x_{\text{com},i} = x_{\text{com},i}^1 + x_{\text{com},i}^2, \quad (27)$$

$$x_{\text{gat},i} = \frac{1}{K} \sum_{k=1}^K x_{\text{com},i}^k \quad (28)$$



$$\begin{aligned}
 v_i &= \text{ReLU}(W_0 h_i + b_0), \\
 p_i &= \text{Softmax}(W_1 v_i + b_1), \\
 \hat{y}_i &= \underset{k}{\text{argmax}}(p_i[k]),
 \end{aligned} \tag{29}$$

$$\mathcal{L}_{cls} = -\frac{1}{\sum_{k=0}^{N-1} n(k)} \sum_{i=0}^{N-1} \sum_{j=0}^{n(i)-1} y_{ij} \log p_{ij} + \lambda |\Theta_{re}|, \tag{30}$$

$$\mathcal{L} = \mathcal{L}_{cls} + \beta \mathcal{L}_{shr} + \gamma^a \mathcal{L}_{sep}^a + \gamma^v \mathcal{L}_{sep}^v + \gamma^t \mathcal{L}_{sep}^t, \tag{31}$$



Experiments

Model	IEMOCAP								MELD						
	<i>Happy</i>	<i>Sad</i>	<i>Neutral</i>	<i>Angry</i>	<i>Excited</i>	<i>Frustrated</i>	Accuracy	wa-F1	<i>Neutral</i>	<i>Surprise</i>	<i>Sadness</i>	<i>Joy</i>	<i>Anger</i>	Accuracy	wa-F1
bc-LSTM	32.63	70.34	51.14	63.44	67.91	61.06	59.58	59.10	75.66	48.47	22.06	52.10	44.39	59.62	56.80
CMN	30.38	62.41	52.39	59.83	60.25	60.69	56.56	56.13	-	-	-	-	-	-	-
ICON	29.91	64.57	57.38	63.04	63.42	60.81	59.09	58.54	-	-	-	-	-	-	-
DialogueRNN	33.18	78.80	59.21	65.28	71.86	58.91	63.40	62.75	76.79	47.69	20.41	50.92	45.52	60.31	57.66
DialogueCRN	51.59	74.54	62.38	67.25	73.96	59.97	65.31	65.34	76.13	46.55	11.43	49.47	44.92	59.66	56.76
DialogueGCN	47.10	80.88	58.71	66.08	70.97	61.21	65.54	65.04	75.97	46.05	19.60	51.20	40.83	58.62	56.36
MMGCN	45.45	77.53	61.99	66.67	72.04	64.12	65.56	65.71	75.16	48.45	25.71	54.41	45.45	59.31	57.82
GraphCFC	43.08	84.99	64.70	71.35	78.86	63.70	69.13	68.91	76.98	49.36	26.89	51.88	47.59	61.42	58.86



Experiments

Modality Setting	IEMOCAP		MELD	
	Accuracy	wa-F1	Accuracy	wa-F1
A	54.16	53.85	47.55	41.62
V	31.61	27.67	47.59	33.26
T	59.95	60.09	60.77	56.81
A + V	54.17	53.89	47.61	41.67
A + T	64.20	64.74	59.96	57.46
V + T	63.15	62.96	59.46	57.29
A + V + T	69.13	68.91	61.42	58.86

Experiments

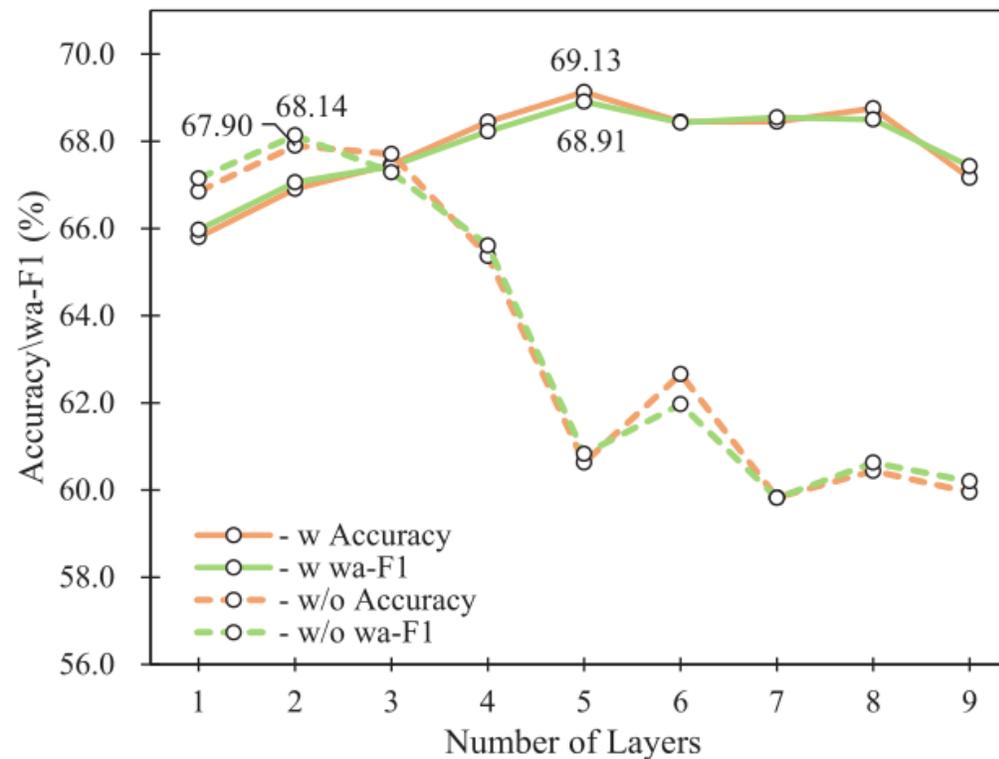


Fig. 4: The effects of the number of GAT-MLP layer and skip connection on our model. The figure shows the results on the IEMOCAP dataset. w (w/o) indicates the use (non-use) of skip connection.



Experiments

TABLE V: The effect of the multi-subspace loss functions. \mathcal{L}_{shr} and $\mathcal{L}_{sep}^{\zeta}$ ($\zeta \in \{a, v, t\}$) denote shared and separate subspace losses, respectively.

\mathcal{L}_{shr}	$\mathcal{L}_{sep}^{\zeta}$	IEMOCAP		MELD	
		Accuracy	wa-F1	Accuracy	wa-F1
- w/o	- w	68.70	68.35	61.00	58.39
- w	- w/o	67.53	67.56	60.27	57.99
- w/o	- w/o	68.70	68.36	60.38	58.09
- w	- w	69.13	68.91	61.42	58.86



Experiments

TABLE VI: The influence of speakers and edge types on our GraphCFC model. S_{emb} and E_{emb} indicate the embeddings of multi-speaker and edge types, respectively.

S_{emb}	E_{emb}	IEMOCAP		MELD	
		Accuracy	wa-F1	Accuracy	wa-F1
- w/o	- w	68.02	68.04	60.69	58.35
- w	- w/o	65.26	65.91	60.46	57.91
- w	- w	69.13	68.91	61.42	58.86

Experiments

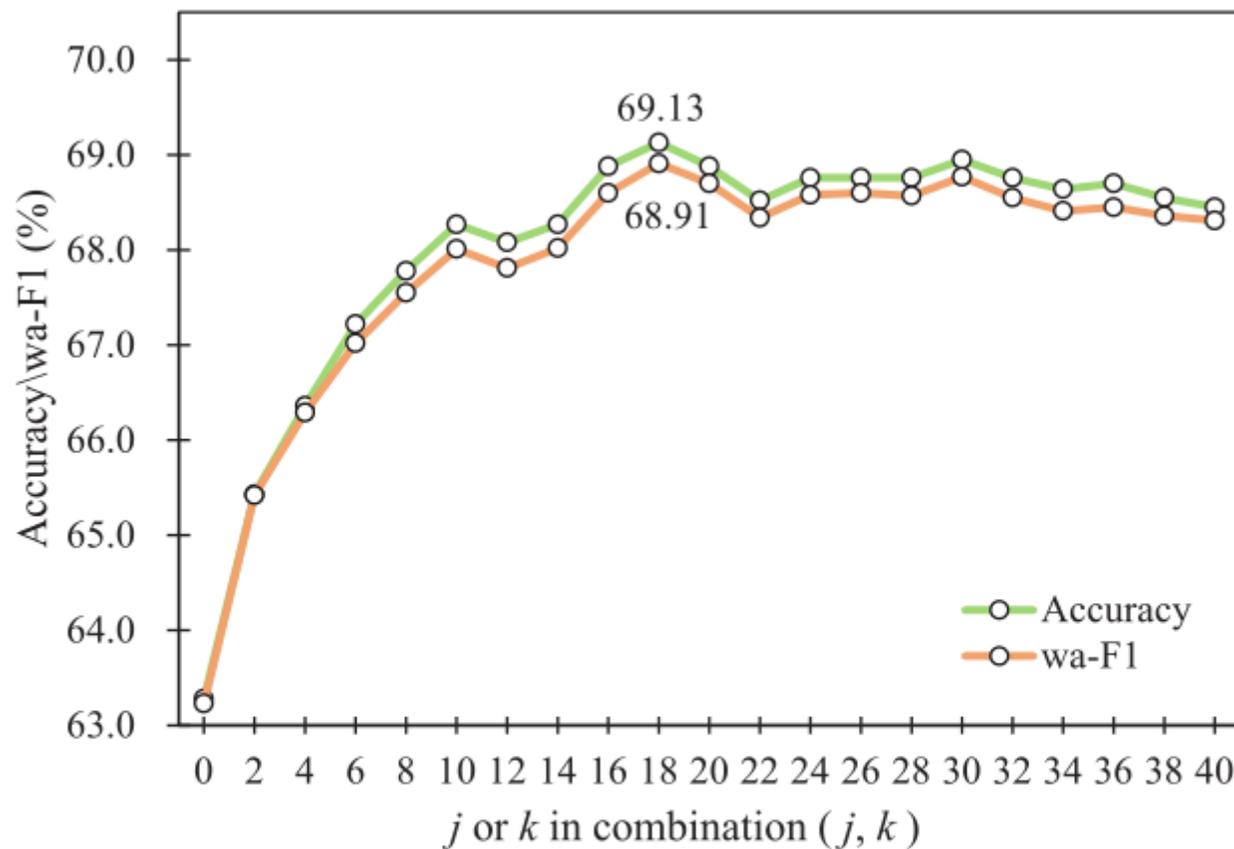


Fig. 5: The effects of j nodes in the past and k nodes in the future on the proposed GraphCFC model. The figure shows the results on the IEMOCAP dataset.



Experiments

TABLE VII: The statistics of the merged emotion labels.

New Label	IEMOCAP	MELD
<i>Positive</i>	<i>Happy, Excited</i>	<i>Joy</i>
<i>Negative</i>	<i>Sad, Angry, Frustrated</i>	<i>Surprise, Fear, Sadness, Disgust, Anger</i>
<i>Neutral</i>	<i>Neutral</i>	<i>Neutral</i>



Experiments

TABLE VIII: The overall performance after converting the dataset into three-emotion labels under the multimodal setting.

Model	IEMOCAP					MELD				
	<i>Positive</i>	<i>Neutral</i>	<i>Negative</i>	Accuracy	wa-F1	<i>Positive</i>	<i>Neutral</i>	<i>Negative</i>	Accuracy	wa-F1
bc-LSTM	90.58	55.63	84.04	79.54	79.10	36.97	75.12	61.46	65.13	64.26
DialogueRNN	88.36	57.99	83.81	78.87	78.94	40.29	74.95	62.10	65.52	64.93
DialogueCRN	79.39	61.51	83.09	75.66	76.97	40.80	74.40	62.87	65.98	65.32
DialogueGCN	84.22	56.88	83.66	77.57	77.48	32.92	75.64	63.96	66.67	64.80
MMGCN	85.20	64.21	83.73	79.36	79.95	43.32	75.5	65.57	67.93	66.92
GraphCFC	88.48	62.03	84.35	79.91	80.20	50.66	75.12	66.26	68.54	68.12

Experiments

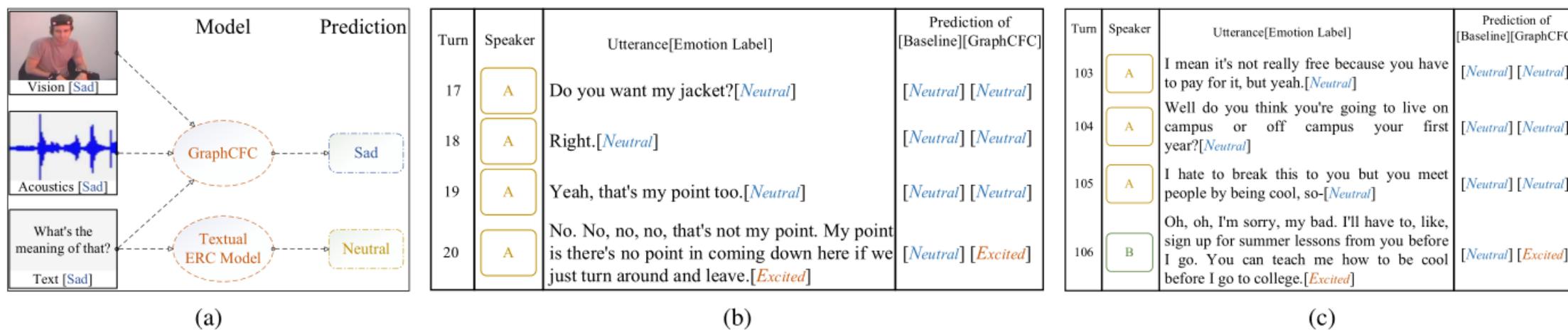


Fig. 6: The cases of ERC on the IEMOCAP. (a) An example shows that multi-modality can be used to compensate for the shortcoming of single-textual modality. (b) Emotional-shift in one-speaker scenario. (c) Emotional-shift in two-speaker scenario.



Thanks!