



Real-time Emotion Pre-Recognition in Conversations with Contrastive Multi-modal Dialogue Pre-training

Xincheng Ju

School of Computer Science and
Technology, Soochow University
Suzhou, Jiangsu, China
xcju@stu.suda.edu.cn

Dong Zhang*

School of Computer Science and
Technology, Soochow University
Suzhou, Jiangsu, China
dzhang@suda.edu.cn

Suyang Zhu

School of Computer Science and
Technology, Soochow University
Suzhou, Jiangsu, China
syzhu@suda.edu.cn

Junhui Li

School of Computer Science and
Technology, Soochow University
Suzhou, Jiangsu, China
lijunhui@suda.edu.cn

Shoushan Li

School of Computer Science and
Technology, Soochow University
Suzhou, Jiangsu, China
lishoushan@suda.edu.cn

Guodong Zhou

School of Computer Science and
Technology, Soochow University
Suzhou, Jiangsu, China
gdzhou@suda.edu.cn

Code:None.

— — CIKM 2023



gesis
Leibniz-Institut
für Sozialwissenschaften



2023.11.23

Reported by Tingting Zhang

Motivation

•SER (Static Emotion Recognition)

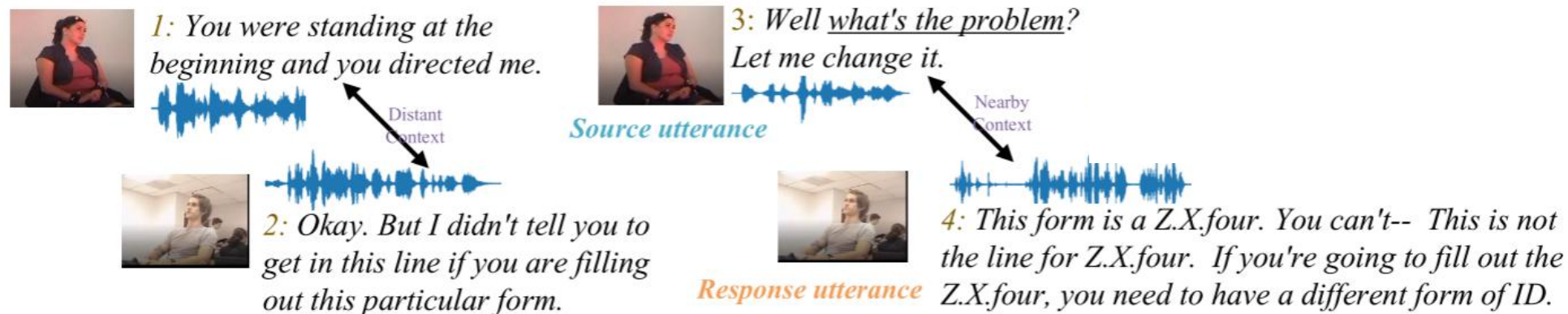
Many studies with focus on static emotion recognition, leverage past and future context to speculate the emotion of the existing target utterance in the textual or multi-modal scenario.

•RERC (Real-time Emotion Recognition in Conversations)

Recently, textual and multi-modal approaches observe the practical applications of RERC, which only leverage the past context to detect the target utterance emotion.

•MREPC (Multi-modal Real-time Emotion Pre-recognition in Conversations)

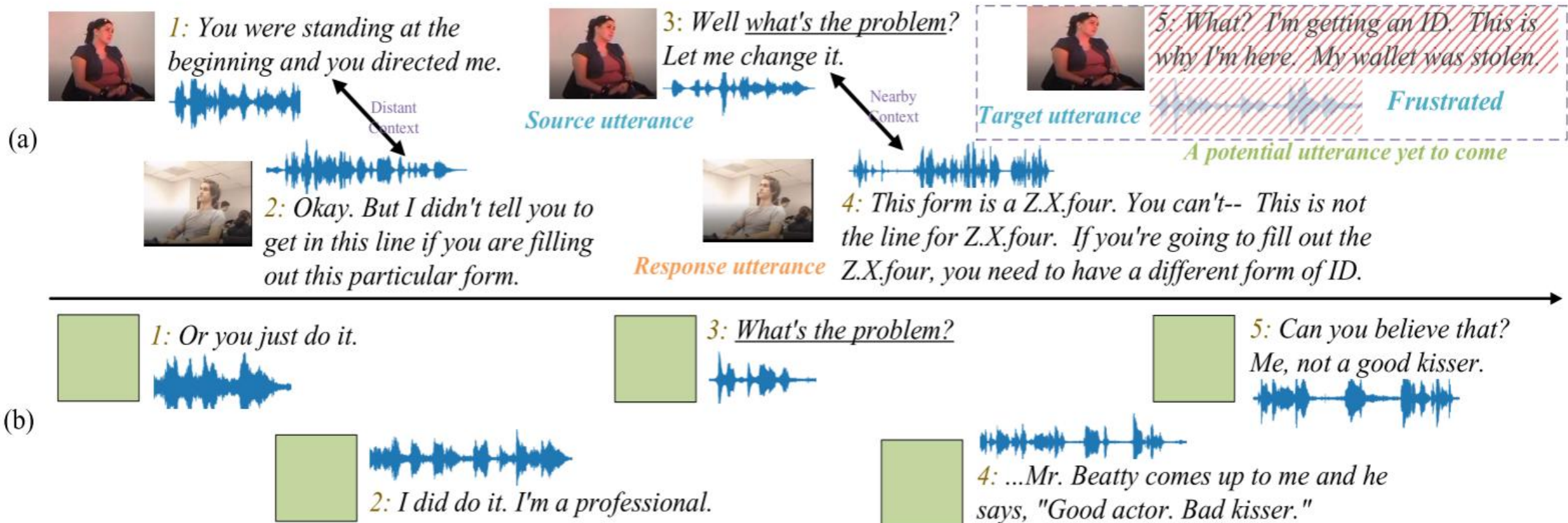
The objective is to predict the emotion of a forthcoming target utterance that is highly likely to occur (lack the target utterance and can only depend on the historical context).



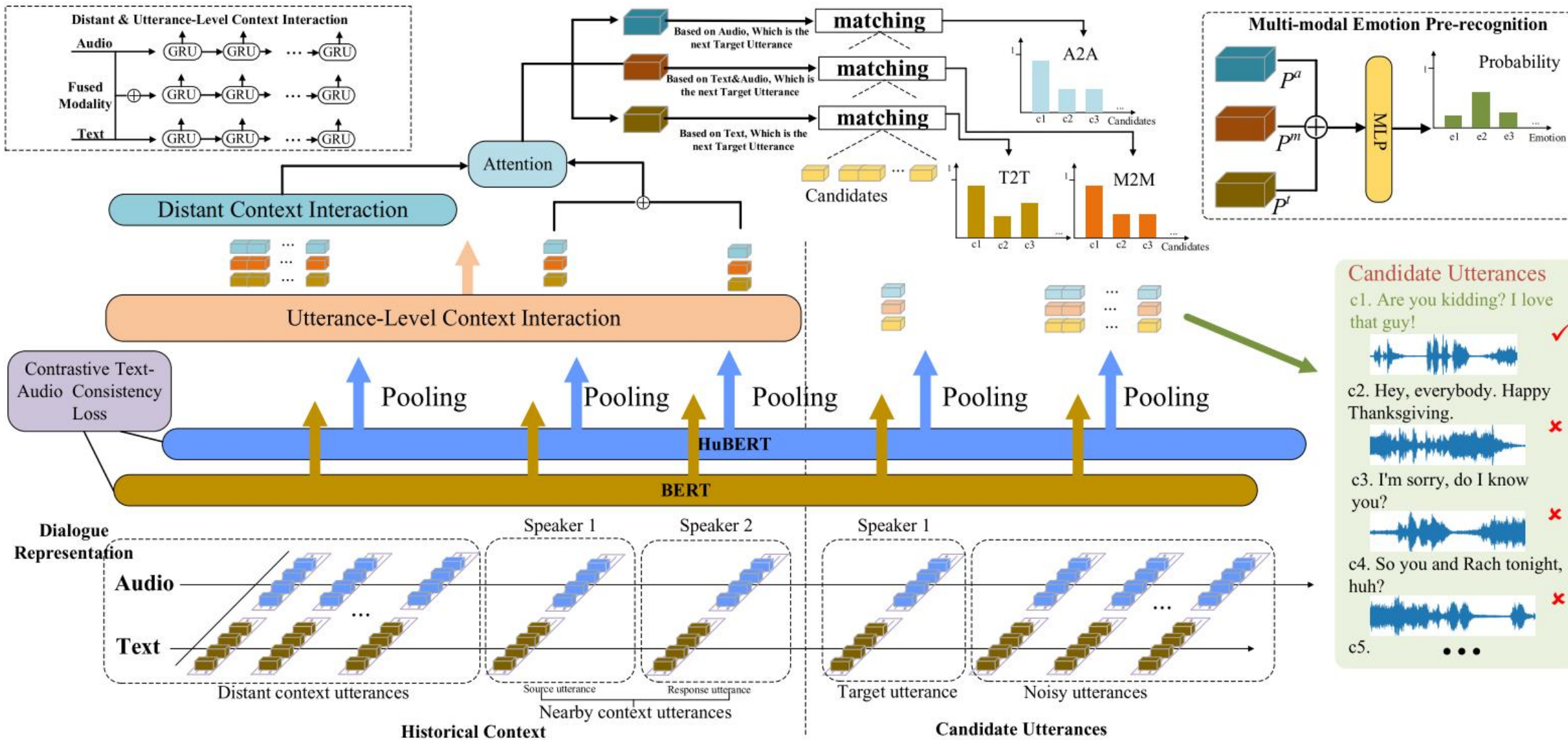
Motivation

It's difficult to speculate the emotion of the target utterance, due to the fact that the target utterance is not existing and the historical context can not provide adequate clues for emotion pre-recognition.

While, a complete conversation though unlabelled in Figure 1(b) may greatly supply the empathic information for (a).



Overview



Method

1. Task Definition

$$p(y|\mathcal{D}) = p(e_{N+1}|C_d, C_n) \quad (1)$$

2. Modality Encoding

$$T_i^t = \text{BERT}(u_i^t) \quad i \in \{1, 2, \dots, N\} \quad (2)$$

$$T_i^a = \text{HuBERT}(u_i^a) \quad i \in \{1, 2, \dots, N\} \quad (3)$$

3. Utterance Representation

$$H_i^{\{t,a\}} = \text{MaxPool}(T_i^{\{t,a\}}) + \text{MeanPool}(T_i^{\{t,a\}}) \quad (4)$$

$$H_c, H_s, H_r = H_{\{1:N-2\}}, H_{N-1}, H_N \quad (5)$$

4. Utterance – level Context Interaction

$$H^m = W_m(H_i^t \oplus H_i^a) + b_m \quad i \in \{1, 2, \dots, N\} \quad (6)$$

$$\hat{H}^\beta = \text{GRU}^\beta(H^\beta) \quad (7)$$

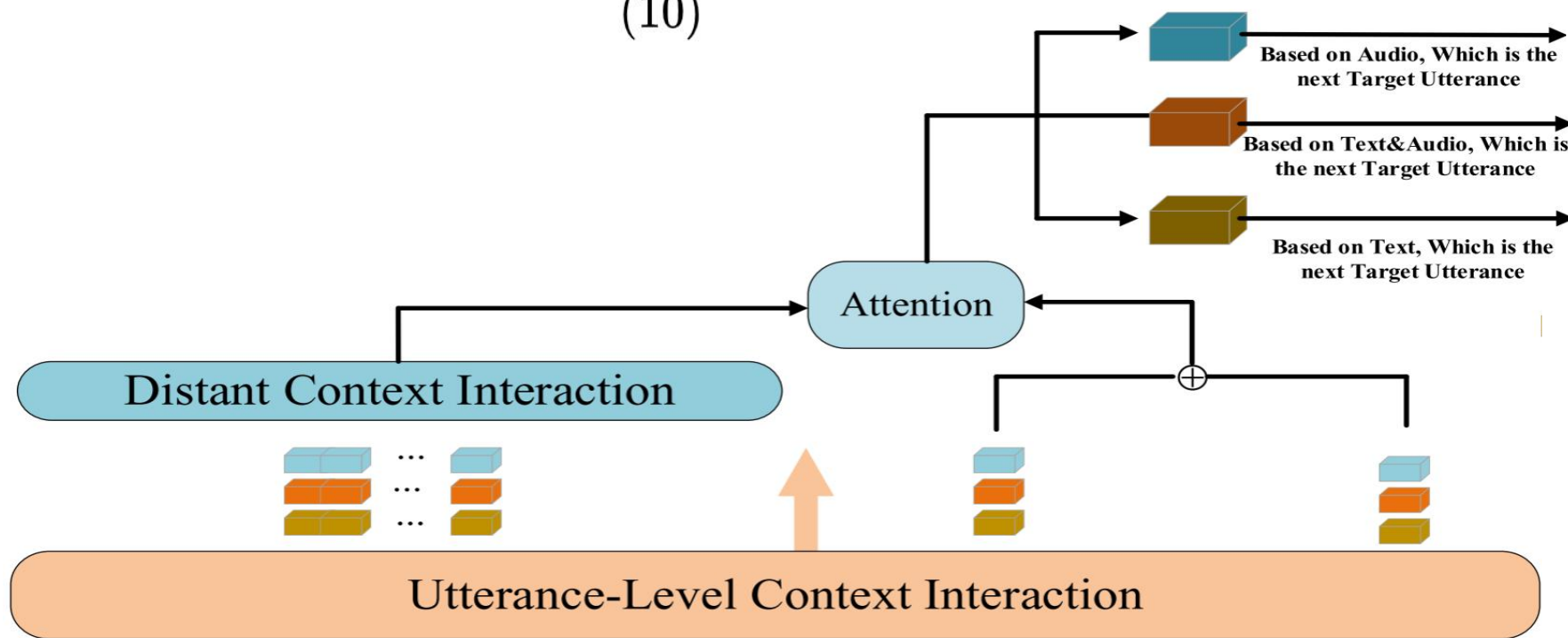
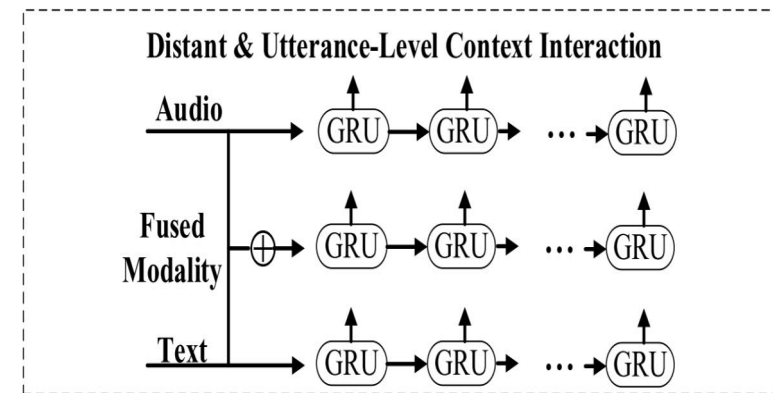
Method

5. Distant and Nearby Context Interaction

$$C^\beta = \text{GRU}^\beta(\hat{H}_c^\beta) \tag{8}$$

$$Q^\beta = (H_s^\beta \oplus H_r^\beta)W_{q\beta} + b_{q\beta} \tag{9}$$

$$P^\beta = \text{Attention}(Q^\beta, C^\beta) \tag{10}$$



Method

1. Intra – modal Target Utterance Searching

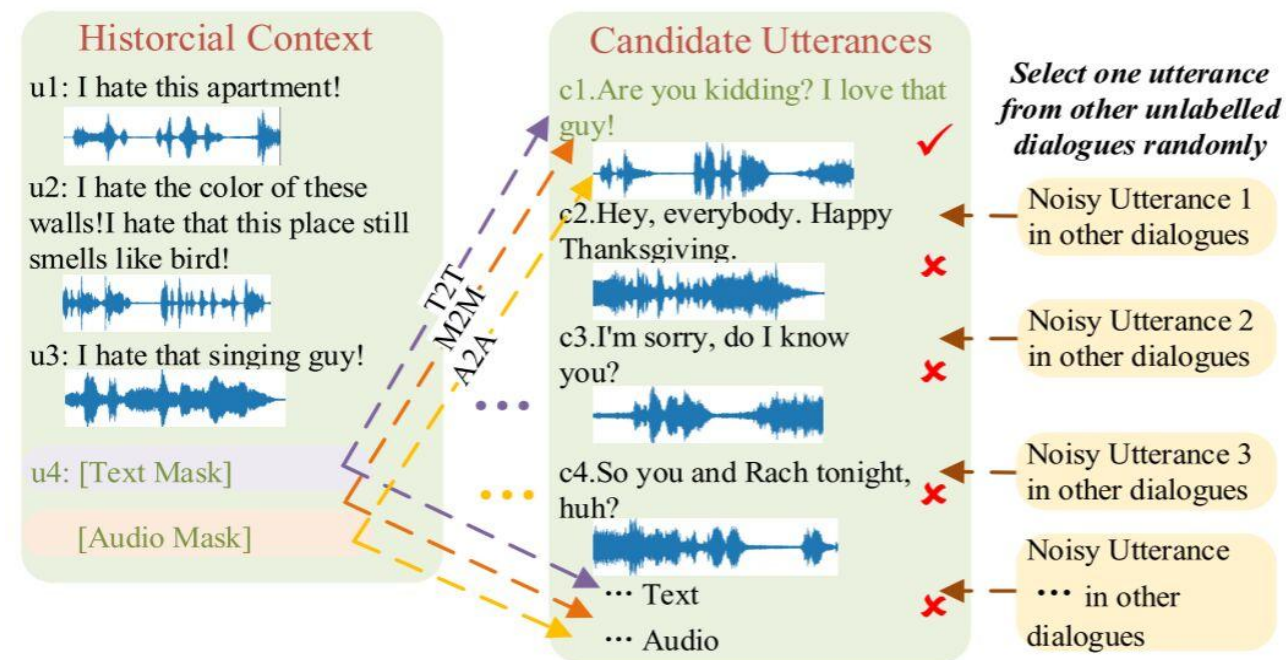
$$s_i^{t \rightarrow t} = \sigma((P^t)^\top K_i^t) \quad i \in \{1, 2, \dots, k\} \quad (11)$$

$$s_i^{a \rightarrow a} = \sigma((P^a)^\top K_i^a) \quad i \in \{1, 2, \dots, k\} \quad (12)$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

2. Inter – modal Target Utterance Searching

$$s_i^{m \rightarrow m} = \sigma((P^m)^\top K_i^m) \quad i \in \{1, 2, \dots, k\} \quad (13)$$





Method

3. Contrastive Loss of Target Utterance Searching

$$\mathcal{L}_{Intra} = -\log(s_i^{t \rightarrow t} \cdot s_i^{a \rightarrow a}) + \sum_{j \in \{1, 2, \dots, k\} - i} \log(s_i^{t \rightarrow t} \cdot s_i^{a \rightarrow a}) \quad (14)$$

$$\mathcal{L}_{Inter} = -\log(s_i^{m \rightarrow m}) + \sum_{j \in \{1, 2, \dots, k\} - i} \log(s_i^{m \rightarrow m}) \quad (15)$$

$$\mathcal{L}_{II} = \mathcal{L}_{Intra} + \mathcal{L}_{Inter} \quad (16)$$

Method

4. Contrastive Loss of Text – Audio Consistency

$$\text{logits} = H^t \cdot (H^a)^\top \tag{17}$$

$$\mathcal{L}_h = - \sum_{i=1}^N \log \left(\frac{\exp^{\text{logits}[i,i]}}{\sum_{j=1}^N \exp^{\text{logits}[i,j]}} \right) \tag{18}$$

$$\mathcal{L}_v = - \sum_{i=1}^N \log \left(\frac{\exp^{\text{logits}[i,i]}}{\sum_{j=1}^N \exp^{\text{logits}[j,i]}} \right) \tag{19}$$

$$\mathcal{L}_{hv} = (\mathcal{L}_h + \mathcal{L}_v) / 2 \tag{20}$$

5. Total Pre – training Loss

$$\mathcal{L}_{total} = \zeta \mathcal{L}_{II} + \eta \mathcal{L}_{hv} \tag{21}$$

| | A ₁ | A ₂ | A ₃ | ... | A _N |
|----------------|---------------------------------|---------------------------------|---------------------------------|-----|---------------------------------|
| T ₁ | T ₁ • A ₁ | T ₁ • A ₂ | T ₁ • A ₃ | ... | T ₁ • A _N |
| T ₂ | T ₂ • A ₁ | T ₂ • A ₂ | T ₂ • A ₃ | ... | T ₂ • A _N |
| T ₃ | T ₃ • A ₁ | T ₃ • A ₂ | T ₃ • A ₃ | ... | T ₃ • A _N |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| T _N | T _N • A ₁ | T _N • A ₂ | T _N • A ₃ | ... | T _N • A _N |

Figure 4: Contrastive Text-Audio Consistency Loss. T:Text, A:Audio

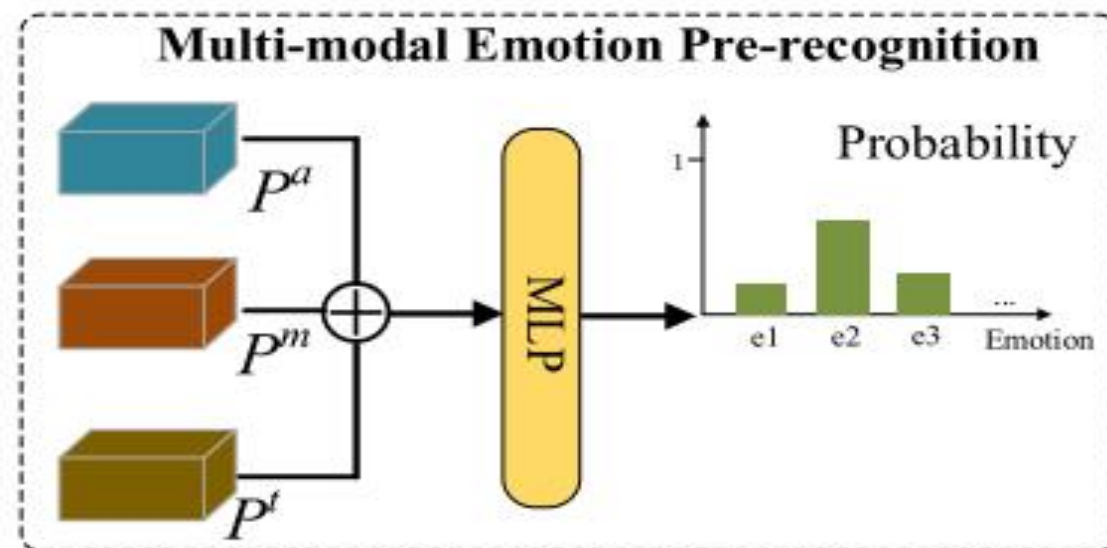
Method

1. Pre – recognition

$$\tilde{U} = \rho(\varrho((P^t \oplus P^m \oplus P^a)W_1 + b_1)W_2 + b_2)) \quad (22)$$

2. Loss of MREPC

$$\mathcal{J} = -\omega(y) \cdot \log(\tilde{U}^y) \quad (23)$$



Experiments

$$\mathcal{L}_{Intra} = -\log(s_i^{t \rightarrow t} \cdot s_i^{a \rightarrow a}) + \sum_{j \in \{1, 2, \dots, k\} - i} \log(s_i^{t \rightarrow t} \cdot s_i^{a \rightarrow a}) \quad (14)$$

$$\mathcal{L}_{Inter} = -\log(s_i^{m \rightarrow m}) + \sum_{j \in \{1, 2, \dots, k\} - i} \log(s_i^{m \rightarrow m}) \quad (15)$$

$$\mathcal{L}_{II} = \mathcal{L}_{Intra} + \mathcal{L}_{Inter} \quad (16)$$

| Approaches | TUS | $R_5@1$ | $R_5@2$ | $R_{11}@1$ | $R_{11}@2$ |
|-------------|-----|---------|---------|------------|------------|
| Intra-modal | T2T | 41.18 | 63.74 | 24.60 | 40.13 |
| | A2A | 96.38 | 97.20 | 96.89 | 96.93 |
| Inter-modal | M2M | 80.82 | 82.79 | 79.96 | 80.91 |
| Multi-modal | T2T | 45.36 | 67.12 | 27.88 | 41.70 |
| | A2A | 87.57 | 89.59 | 86.05 | 87.50 |
| | M2M | 90.59 | 92.84 | 89.33 | 90.46 |

Table 3: The performance of target utterance searching (TUS) in pre-training with different perspectives. The intra-modal part comes from our model which only calculates the intra-modal TUS loss (Eq. 14), while the inter-modal part only calculates the inter-modal TUS loss (Eq. 15). Multi-modal part denotes that both the intra- and inter-modal TUS loss are calculated simultaneously (Eq. 16). Note that higher scores mean better experimental results.

Experiments

| Modality | Approaches | P-MELD | | | P-IEM | | |
|-------------------|--------------------|-------------|-------------|-------------|-------------|-------------|------|
| | | WA | AA | F1 | WA | AA | F1 |
| Text | PPE-Text | 45.5 | 33.3 | 20.9 | 29.4 | 10.0 | 4.5 |
| | NSF-Text-our impl. | 48.6 | 41.3 | 35.6 | 45.2 | 17.9 | 12.0 |
| | AGHMN-Text | 43.4 | 36.4 | 32.1 | 41.7 | 20.7 | 19.4 |
| | DialogXL-Text | 45.5 | 33.3 | 20.9 | 44.9 | 20.5 | 18.8 |
| Text+Audio | PPE[9] | 43.3 | 33.9 | 29.2 | 34.8 | 12.6 | 9.3 |
| | NSF-our impl. [39] | 49.2 | 41.3 | 36.4 | 44.4 | 17.5 | 11.7 |
| | AGHMN [15] | 44.7 | 33.7 | 30.5 | 35.4 | 13.1 | 9.5 |
| | DialogXL [35] | 42.4 | 35.3 | 31.3 | 45.3 | 21.1 | 18.3 |
| | DialogXL-our impl. | 46.3 | 36.9 | 32.5 | 48.8 | 27.9 | 26.4 |
| | BiDDIN [42] | 42.2 | 34.4 | 30.5 | 34.3 | 14.7 | 12.7 |
| | MMGCN [13] | 46.5 | 34.6 | 24.7 | 36.9 | 22.8 | 20.5 |
| | MDI-our impl. [44] | 46.7 | 40.0 | 35.2 | 42.6 | 17.3 | 11.8 |
| TCMP(ours) | 50.1 | 41.5 | 38.9 | 53.9 | 28.2 | 27.4 | |

Table 4: The performance of different approaches for MREPC task. Text: only utilize textual modality. our impl.: implementing our modified setting for the corresponding approach.



Experiments
















| Approaches | P-MELD | | | P-IEM | | |
|---------------------|-------------|------|-------------|-------------|-------------|-------------|
| | WA | AA | F1 | WA | AA | F1 |
| Text w/o Dist | 49.2 | 41.2 | 36.7 | 44.8 | 17.8 | 11.9 |
| Text | 49.5 | 41.7 | 36.9 | 44.6 | 17.6 | 11.8 |
| Text-Pre | 49.3 | 41.1 | 37.8 | 48.9 | 21.2 | 18.8 |
| Text+Audio w/o Dist | 48.4 | 41.5 | 37.1 | 46.0 | 19.6 | 15.8 |
| Text+Audio | 49.6 | 42.1 | 36.8 | 52.7 | 27.3 | 26.8 |
| TCMP(ours) | 50.1 | 41.5 | 38.9 | 53.9 | 28.2 | 27.4 |

Table 5: The performance of single-modal and multi-modal ablated approaches on both datasets.

-Pre denotes the pre-trained approach.

w/o Dist denotes the elimination of distant context.

Experiments

| | Golden | Frustrated | Sad | Excited |
|------------|--|---|--|--|
| Dialogues | <p>1: I went to school and I got my degree and I got a job</p>  <p>2: I mean I just do not know if you do not have a lot of qualification, where do you find work? it's so frustrating because if you don't know somebody you can not get a job, it's totally discriminatory you have to know somebody</p>  <p>3: nothing is impossible</p>  | <p>4: it's not fair</p>  | <p>1: What's he going to say? Maybe we should tell him before he sees it.</p>  <p>2: He saw it.</p>  <p>3: How could he see it? I was the first one up. He was still in bed.</p>  <p>4: he was out here when it broke</p>  | <p>1: Were you in Chicago?</p>  <p>2: Mm-hmm.</p>  <p>3: Downtown, was it beautiful? of course it was beautiful</p>  <p>4: Yeah. so beautiful, full moon.</p>  <p>5: What a guy. I can't believe it. So okay, so do you know any details? When's it going to be? Anything?</p>  <p>6: I don't know. I guess next summer.</p>  |
| | | | <p>5: when?</p>  | |
| AGHMN | | Angry | Neutral | Neutral |
| MMGCN | | Frustrated | Surprise | Neutral |
| Text+Audio | | Frustrated | Sad | Happy |
| TMCP | | Frustrated | Sad | Excited |
| | | (a) | (b) | (c) |



Thanks!