



# Speaker-Guided Encoder-Decoder Framework for Emotion Recognition in Conversation

Yinan Bao<sup>1,2</sup>, Qianwen Ma<sup>1,2</sup>, Lingwei Wei<sup>1,2</sup>, Wei Zhou<sup>1,\*</sup> and Songlin Hu<sup>1,2</sup>

<sup>1</sup>Institute of Information Engineering, Chinese Academy of Sciences

<sup>2</sup>School of Cyber Security, University of Chinese Academy of Sciences

{baoyinan, maqianwen, weilingwei, zhouwei, husonglin}@iie.ac.cn

2022. 6. 29 • ChongQing

— IJCAI 2022



gesis  
Leibniz-Institut  
für Sozialwissenschaften



Reported by Sijin Liu



# 1. Introduction

## 2. Method

### 3. Experiments



# Introduction

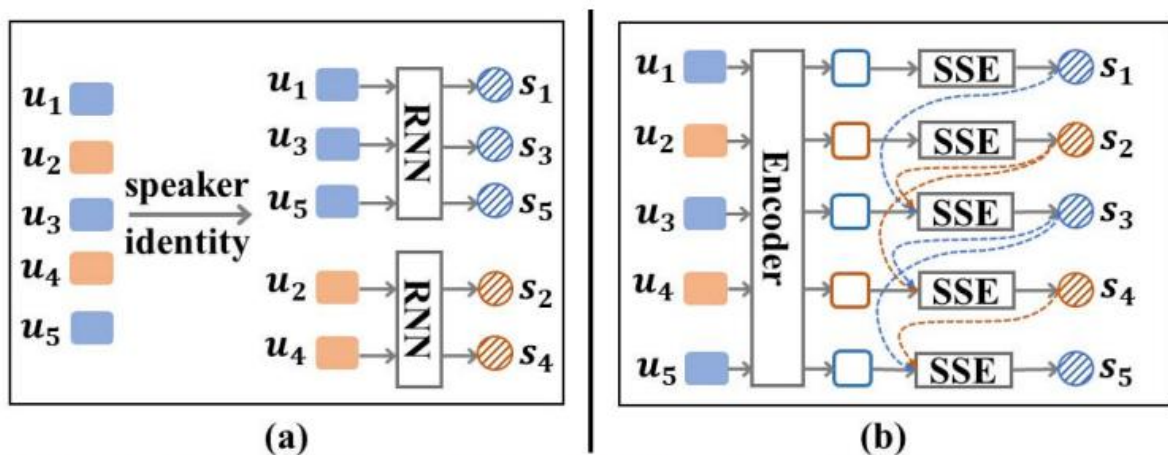


Figure 1: (a) Speaker modeling scheme of existing dynamic speaker-specific modeling methods, which is based on speaker identity to explore intra-speaker dependency. (b) Our novel speaker modeling scheme, which uses the speaker state encoder (SSE) to explore intra- and inter-speaker dependencies dynamically.  $u_i$  means the  $i_{th}$  utterance in a dialogue and  $s_i$  means the corresponding speaker state vector of  $u_i$ .

Sequences composed of separate utterances attaching to the same speaker are **disjunctive and incomplete**, which hinder the comprehension of the context.

**Overreliance on the speaker identity** hinders the modeling of **dynamic inter-speaker dependency**. In this way, the model can't explore dynamic interactions between different speakers.

# Method

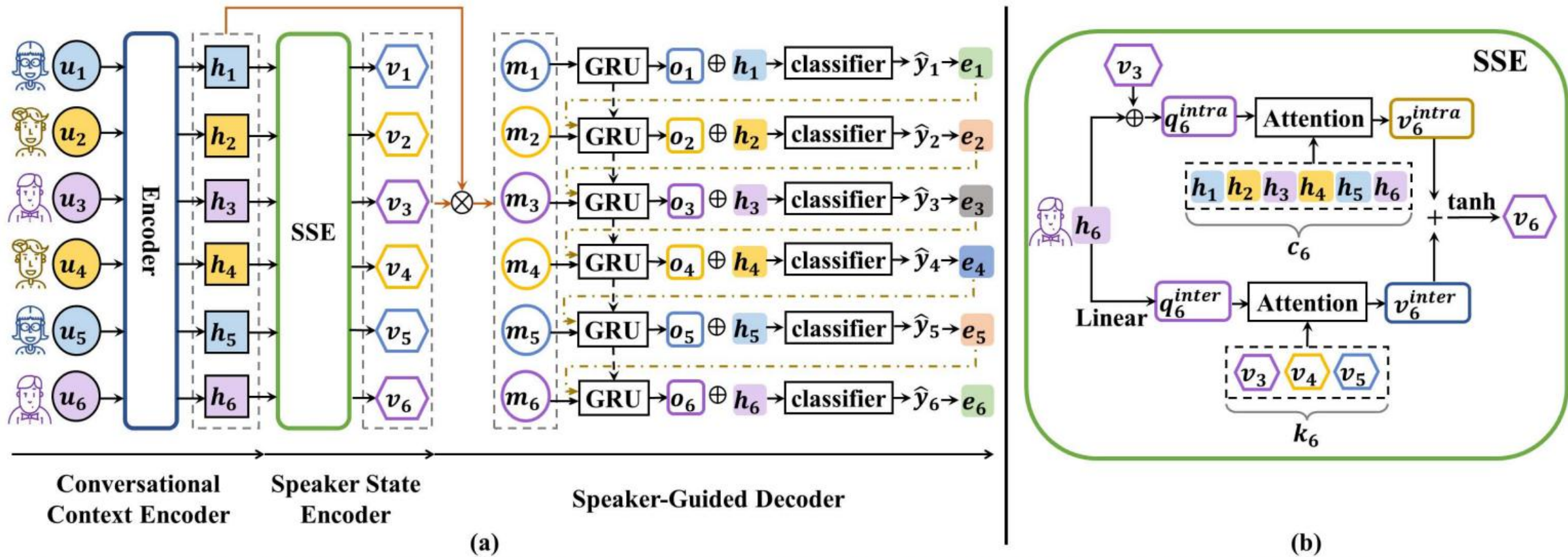
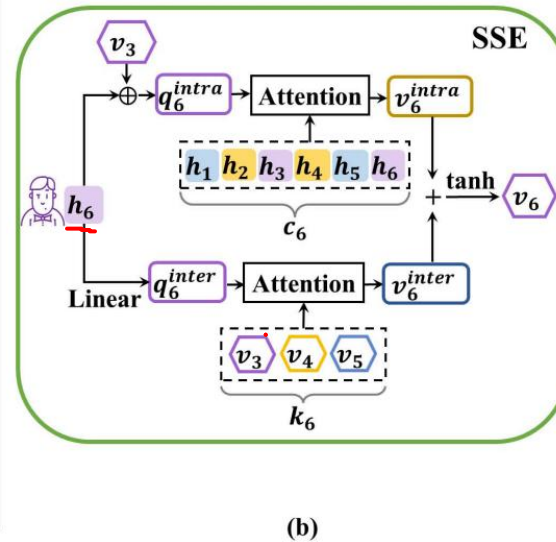
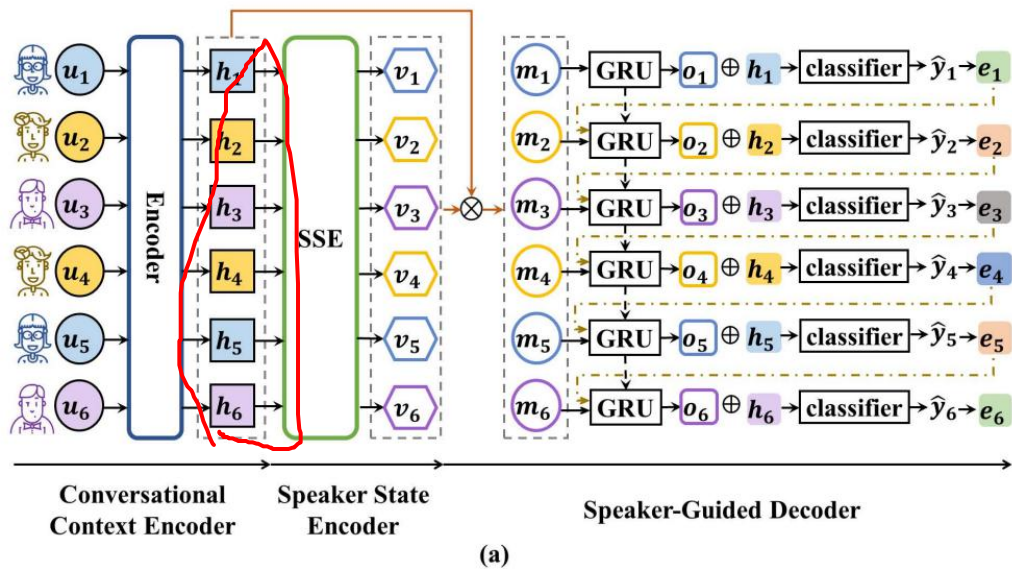


Figure 2: An overview of SGED. SSE represents the speaker state encoder. As shown in the right part of the figure, we take  $u_6$  as an example to show how to obtain the speaker state vector of it.

# Method



## Conversational Context Encoder

$$C = [u_1, \dots, u_N]$$

$$H = \text{Encoder}(C), \quad (1)$$

## Speaker State Encoder

### Intra-Speaker Dependency Modeling

$$q_i^{intra} = W_q^{intra} [v_{\psi(u_i)} || h_i] + b_q^{intra}, \quad (2)$$

where  $v_{\psi(u_i)} \in \mathbb{R}^{h \times 1}$  means the last previous speaker state vector of  $s_{\phi(u_i)}$ ;  $||$  means the concatenation operation;  $W_q^{intra} \in \mathbb{R}^{h \times 2h}$  and  $b_q^{intra} \in \mathbb{R}^1$  are model parameters.

$$\alpha_{ij}^{intra} = \text{softmax}(W_1(q_i^{intra} \odot c_j) + b_1), \quad (3)$$

$$v_i^{intra} = \sum_j \alpha_{ij}^{intra} c_j^T,$$

### Inter-Speaker Dependency Modeling

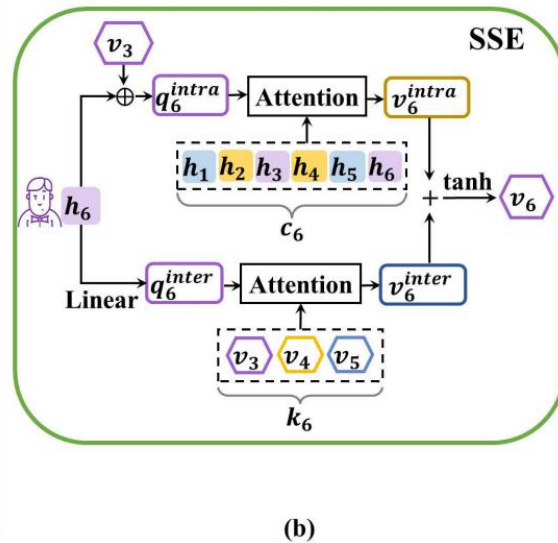
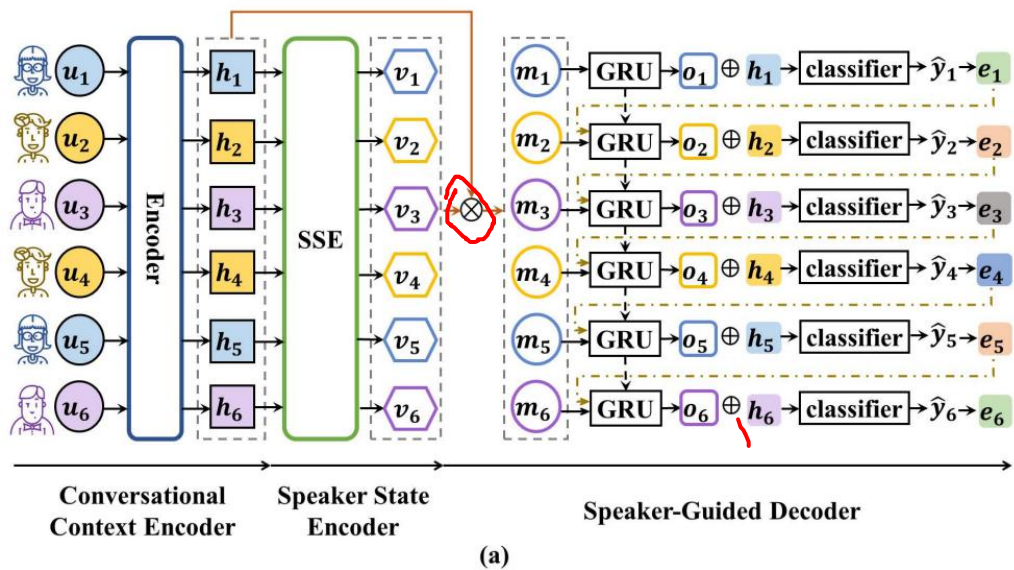
$$q_i^{inter} = W_q^{inter} h_i + b_q^{inter},$$

$$\alpha_i^{inter} = \text{softmax}(W_2(q_i^{inter} \odot k_i) + b_2), \quad (4)$$

$$v_i^{inter} = \alpha_i^{inter} k_i^T,$$

$$v_i = \tanh(v_i^{intra} + v_i^{inter}), \quad (5)$$

# Method



## Speaker-Guided Decoder

$$\mathbf{m}_i = \text{ReLU}(\mathbf{v}_i \odot (\mathbf{W}_m \mathbf{h}_i + \mathbf{b}_m)^T), \quad (6)$$

$$\mathbf{o}_i = \overrightarrow{\text{GRU}}([\mathbf{m}_i || \mathbf{e}_{i-1}], \mathbf{o}_{i-1}). \quad (7)$$

$$\mathbf{z}_i = \text{ReLU}(\mathbf{W}_o [\mathbf{h}_i || \mathbf{o}_i^T] + \mathbf{b}_o),$$

$$\mathcal{P}_i = \text{softmax}(\mathbf{W}_z \mathbf{z}_i + \mathbf{b}_z),$$

$$\hat{y}_i = \underset{g}{\text{argmax}} \mathcal{P}_i[g], \quad (8)$$

$$\mathbf{e}_i = \mathbf{E}[\hat{y}_i],$$

$$\mathcal{L}(\theta) = - \sum_{j=1}^M \sum_{t=1}^{N_j} \log \mathcal{P}_{j,t}[y_{j,t}], \quad (9)$$

# Experiments

Dataset	Conversations			Utterances		
	Train	Val	Test	Train	Val	Test
IEMOCAP	120		31	5810		1623
MELD	1038	114	280	9989	1109	2610
EmoryNLP	713	99	85	9934	1344	1328

Table 1: The statistics of three datasets.

Model	MELD	EmoryNLP	IEMOCAP
KET	58.18	33.95	59.56
COSMIC	65.21	38.11	65.28
DialogueRNN + RoBERTa	57.03 63.61	- 37.44	62.75 64.76
DialogueCRN + RoBERTa	58.39 63.42	- 38.91	66.20 66.46
RoBERTa	62.88	37.78	63.38
<b>SGED + RoBERTa</b>	<b>63.34</b>	<b>38.47</b>	<b>64.11</b>
bc-LSTM+att + RoBERTa	- 62.95	- 38.28	- 64.51
<b>SGED + bc-LSTM+att</b>	<b>63.37</b>	<b>38.89</b>	<b>65.03</b>
DialogueGCN + RoBERTa	58.10 63.02	- 38.10	64.18 64.91
<b>SGED + DialogueGCN</b>	<b>64.55</b>	<b>39.73</b>	<b>65.90</b>
DAG-ERC	63.65	39.02	68.03
DAG-ERC*	63.39	38.84	67.45
<b>SGED + DAG-ERC*</b>	<b>65.46</b>	<b>40.24</b>	<b>68.53</b>

Table 2: Overall performance on the three datasets. We choose weighted-average F1 to evaluate each method. DAG-ERC\* means that we use the one-layer DAG-ERC.



# Experiments

Method	EmoryNLP	IEMOCAP
<b>SGED + DAG-ERC*</b>	<b>40.24</b>	<b>68.53</b>
w/o <b>SSE</b> - intra-speaker	39.87 (↓0.37)	68.05 (↓0.48)
w/o <b>SSE</b> - inter-speaker	39.86 (↓0.38)	67.72 (↓0.81)
w/o <b>SSE</b>	39.17 (↓1.07)	67.67 (↓0.86)
w/o <b>SGD</b>	38.96 (↓1.28)	67.30 (↓1.23)
w/o <b>SSE + SGD</b>	38.84 (↓1.40)	67.45 (↓1.08)

Table 3: Results of ablation study on EmoryNLP (multi-party) and IEMOCAP (two-party) dataset. We use the one-layer DAG-ERC as the conversational context encoder of our framework.



# Experiments

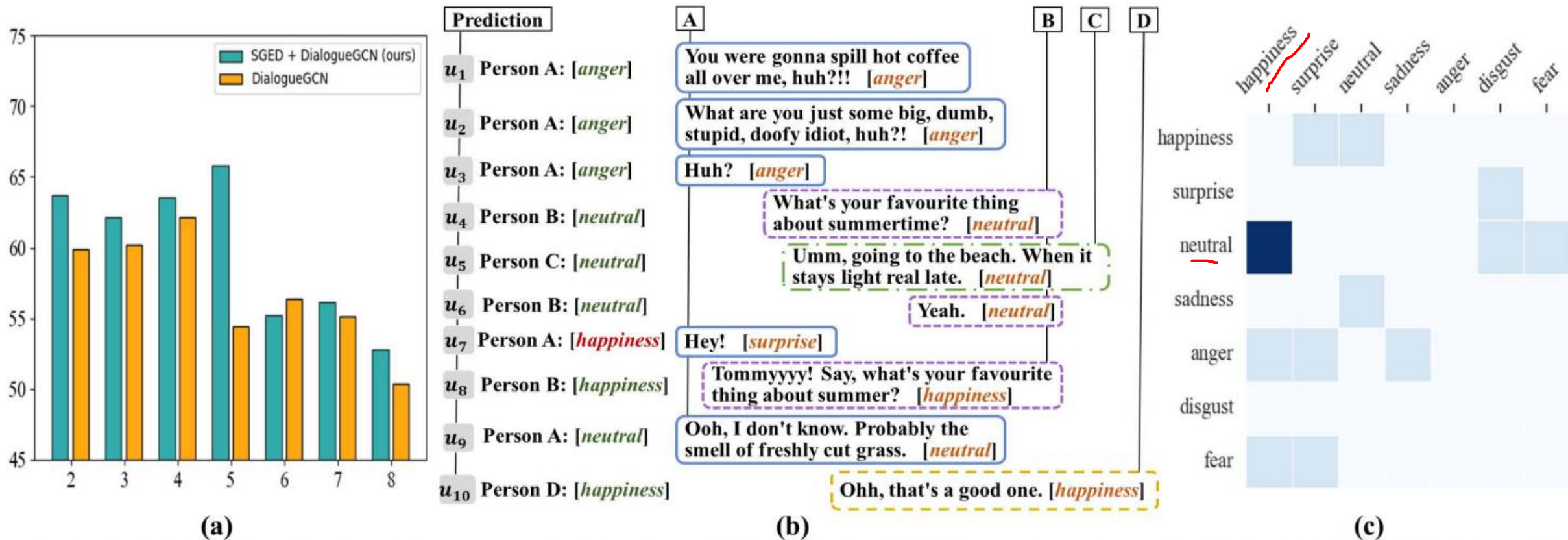


Figure 3: (a) The weighted-average F1 score of conversations with the different number of speakers on MELD, achieved by DialogueGCN and our SGED framework which takes DialogueGCN as the conversational context encoder. (b) An abridged dialogue from MELD in which our SGED framework recognizes most of the emotions correctly. (c) The confusion matrix of conversations with 6 speakers excluding the diagonal on MELD, is achieved by our SGED framework which takes DialogueGCN as the conversational context encoder.