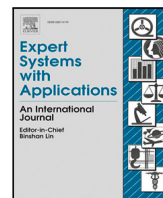




Contents lists available at ScienceDirect

## Expert Systems With Applications

journal homepage: [www.elsevier.com/locate/eswa](http://www.elsevier.com/locate/eswa)

## Knowing What and Why: Causal emotion entailment for emotion recognition in conversations

Hao Liu<sup>a</sup>, Runguo Wei<sup>b</sup>, Geng Tu<sup>b</sup>, Jiali Lin<sup>c</sup>, Dazhi Jiang<sup>b</sup>, Erik Cambria<sup>d</sup>

<sup>a</sup> School of Computer Science, Northwestern Polytechnical University, Xi'An, China

<sup>b</sup> Department of Computer Science of Shantou University, Shantou, China

<sup>c</sup> Business School of Shantou University, Shantou, China

<sup>d</sup> College of Computing and Data Science, Nanyang Technological University, Singapore

### ARTICLE INFO

#### Keywords:

Emotion recognition  
Causal emotion entailment  
Reasoning network

### ABSTRACT

The clues for eliciting emotion deserve attention in the realm of Emotion Recognition in Conversations (ERC). In an ideal dialog system, comprehending emotions alone is insufficient, and underlying the causes of emotion is also imperative. However, previous research overlooked the integration of causal emotion entailment for a prolonged period. Therefore, an emotion-cause hybrid framework that utilizes causal emotion entailment (CEE) is proposed to promote the ERC task. Specifically, the presented method integrates the information of the cause clause extracted through the CEE module that triggers emotions into the utterance representations obtained by the ERC model. Moreover, a Bidirectional Reasoning Network (BRN) is designed to extract emotional cues to simulate human complex emotional cognitive behavior. Experimental results demonstrate that our framework achieves a new state-of-the-art performance on different datasets, indicating that the proposed framework can improve the model's ability to emotion understanding.

### 1. Introduction

Emotion recognition in conversations (ERC) plays a pivotal role in the field of Artificial Intelligence (AI) (Cambria, Howard, Hsu, & Hussain, 2013; Cambria, Mao, Chen, Wang, & Ho, 2023; Koolagudi & Rao, 2012). For example, ERC can be implemented in human-computer interaction, opinion mining, sarcasm detection, etc. Cambria, Zhang, Mao, Chen, and Kwok (2024), Hazarika, Poria, Gorantla et al. (2018), Lee and Hong (2016), Liu, Wei et al. (2024), Liu, Yang and Yu (2024) and Zhu, Mao, Cambria, and Jansen (2024). The emotional content of an utterance is influenced by various factors, such as the conversational context and Causal Emotion Entailment (CEE) (Majumder et al., 2019; Poria et al., 2021). Existing research on ERC mainly uses recurrent neural networks (Hazarika, Poria, Mihalcea, Cambria and Zimmermann, 2018; Majumder et al., 2019) to obtain the dependencies between utterances or use graph-based structure (Saxena, Huang, & Kurohashi, 2022; Shen, Wu, Yang, & Quan, 2021) to gain long-term information. Additionally, transformer-based models are also employed in this task (Chudasama et al., 2022; Luo et al., 2024; Tu, Niu, Xu, Liang, & Cambria, 2024). However, these methods neglect to uncover the causes of emotion generation and its utilization, failing to understand and utilize the information associated with emotions entirely.

Fig. 1 illustrates the connection between cause clauses and emotion clauses. Moreover, these methods tend to ignore partial context information when extracting contextual utterance representation. To address the above issues, we propose a new framework named CauERC with introducing two modules, namely Causal Emotion Entailment (CEE) and Bidirectional Reasoning Network (BRN), into the ERC model.

These modules extract causal information triggering emotions in utterances, offering a significant opportunity to address the identified issues. Recent works for CEE tasks are based on graph networks. Poria et al. (2021) set some baselines for CEE tasks, such as ECPE-MLL, and RankCP, which all use graph attention networks to extract relations between utterances. Zhang, Yang, Meng, Chen, and Zhou (2022) used graph neural networks to provide interaction between utterances and integrate speaker information.

In addition, some works, such as KEC and KBCIN (Li, Meng et al., 2022; Zhao, Zhao, Li, & Qin, 2023) introduce commonsense knowledge into graph neural networks to improve the information extraction ability of the model. MPEG (Chen, Shen, Chen, Zhang, & Zhao, 2023) fuses speaker and sentiment information via a heterogeneous graph attention network to capture the inter-utterances causal relationship. Unlike the methods mentioned above, this paper employs a window transformer.

\* Corresponding authors.

E-mail addresses: [liuhao23@mail.nwpu.edu.cn](mailto:liuhao23@mail.nwpu.edu.cn) (H. Liu), [dzjiang@stu.edu.cn](mailto:dzjiang@stu.edu.cn) (D. Jiang), [cambria@ntu.edu.sg](mailto:cambria@ntu.edu.sg) (E. Cambria).

<sup>1</sup> The authors contribute equally to this work.

<https://doi.org/10.1016/j.eswa.2025.126924>

Received 26 March 2024; Received in revised form 3 February 2025; Accepted 14 February 2025

Available online 22 February 2025

0957-4174/© 2025 Elsevier Ltd. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

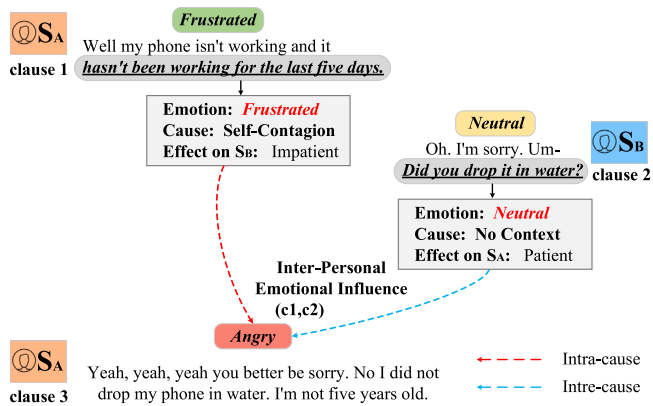


Fig. 1. Example of causal emotion entailment. The dotted line indicates that the emotion of the specified utterance is influenced by the cause clause associated with it. The first utterance of Speaker A indicates that he has already been impatient with Speaker B. But Speaker B apologizes and misinterprets Speaker A's question, which makes Speaker A annoyed again. Because the emotion-cause pairs between utterances like Speaker A's first utterance and Speaker B's second utterance directly triggered Speaker A's anger. This contextual pairing of emotion-cause helps predict participants' emotion labels.

Transformers are particularly effective at capturing the contextual relationships within conversations. Moreover, causal clauses often occur in proximity to emotional clauses, and the window transformer can efficiently model the interaction between inter-utterance semantic information within a limited window size. This makes it well-suited for understanding the nuanced connections between different parts of the conversation. Specifically, in CEE, the utterance is first encoded through the encoding layer. Then we model the inter-clause document through the 2D Window Transformer (Ding, Xia, & Yu, 2020), which is proficient in effectively extracting the semantic correlation between emotion and their underlying cause clauses in conversations. Additionally, inspired by affective cognitive theory, which suggests that individuals infer the emotional states of others by observing their emotional expressions (Lerner, Li, Valdesolo, & Kassam, 2015; Pessoa, 2008), we conceptualize the cognitive process in conversation as the interaction of utterance information between speakers.

By reasoning about the context of these utterances, we can learn the emotional influence between them, simulating the cognitive behaviors that occur in human conversation. To dynamically simulate human emotional and cognitive behavior, we utilize LSTM networks to capture the conversational context and integrate it into bidirectional reasoning. Initially, the contextual information at different stages of the emotion analysis process is stored in static memory nodes. Using LSTM, we integrate and extract this contextual information to understand the underlying logic of conversational utterances and identify emotional cues. At the same time, we update the memory information dynamically to reflect the evolving nature of the conversation. Through multiple iterative reasoning processes, we perform conscious emotional and cognitive reasoning, thereby enhancing the model's ability to accurately simulate human emotional and cognitive behavior in dialog. To summarize, this paper makes the following contributions:

1. We first combine the CEE module with the ERC task so that the model can use the information of the cause clause associated with the utterance used for emotion prediction.
2. We propose the BRN module to imitate human emotion and cognitive behavior in dynamic conversations.
3. Experiments on different conversational datasets showcase that our proposed approach enhances multiple baselines and surpasses state-of-the-art ERC methods.

## 2. Related work

Deep learning plays a significant role in human activities (Fan, Lin, Mao, & Cambria, 2024; Huang, Liu, Tsang, Xu, & Lv, 2022; Jia, Lv, Peng, Chen, & Yang, 2020; Wang, Yu et al., 2023). Emotion analysis around conversations is an important topic in recent years, which has attracted much attention in natural language processing. The availability of many conversation datasets partly explains this phenomenon, and the growing interest in conversational emotion-cause pairs can also explain this phenomenon. In the following paragraphs, we divide the related works into two categories according to the problems they use to model the conversation context.

### 2.1. Emotion recognition in conversations

Rosalind (Picard, 2010) proposes that emotion analysis is an interdisciplinary science that involves psychology, cognitive science, and deep learning. Erik et al. has conducted a comprehensive and proactive exploration of emotional analysis (Susanto, Livingstone, Ng, & Cambria, 2020). With the widespread use of convolutional neural networks (Sun, Yen, Xue, Zhang, & Lv, 2021) and generative adversarial networks (Liu, Fu, Qu, & Lv, 2019; Tang, He, Li, & Lv, 2021), deep learning is also applied in affective computing. Considering the dynamic interaction between speakers, some researchers (Hazarika, Poria, Mihalcea et al., 2018; Majumder et al., 2019) leverage a recurrent neural network to model different speakers to obtain context information. Jiang, Liu, Wei and Tu (2023) applied fuzzy neural network to emotion detection. Due to the recurrent neural network having a long-term information propagation problem, DialogueGCN (Ghosal, Majumder, Poria, Chhaya, & Gelbukh, 2019), DAG-ERC (Shen et al., 2021) and HSGCF (Wang, Dong et al., 2023) employ graph convolution neural network and directed acyclic graph to model the dialog context and simulate the information interaction between speakers, respectively. DualGATs (Zhang, Chen, & Chen, 2023) constructs a dual graph network.

To enrich the utterance representation, KET (Zhong, Wang, & Miao, 2019), SKSEC (Tu, Liang, Jiang and Xu, 2023), CKCL (Tu, Liang, Mao, Yang and Xu, 2023) and COSMIC (Ghosal, Majumder, Gelbukh, Mihalcea, & Poria, 2020) introduce external knowledge into the emotion analysis model by using Knowledge Graph, such as ConceptNet (Liu & Singh, 2004) and COMET (Bosselut et al., 2019), while TODKAT (Zhu, Pergola, Gui, Zhou, & He, 2021) carries out topic detection, and integrates commonsense into a transformer to obtain richer context representation. Li, Zhu, Mao, and Cambria (2023) proposed a Knowledge Integrated Model. To alleviate the issue of category imbalance in emotional data, Tu, Jing et al. (2023) introduced label bias. However, they cannot deal with the problems of difficulty in distinguishing similar emotions and emotion transfer. Therefore, Yang, Shen, Mao, and Cai (2022) constructed a hybrid learning architecture to alleviate the problems of emotion transfer and confusion labeling in conversational emotion. SACL (Hu, Bao, Wei, Zhou, & Hu, 2023) propose the Supervised Adversarial Contrastive Learning to learn structured representations between classes. Multitask learning (Jiang et al., 2021; Tu et al., 2022), self-supervised learning (Jiang, Liu, Tu, Wei, & Cambria, 2024) and contrastive learning (Tu, Liang, Mao et al., 2023) are also applied in emotion identification.

### 2.2. Causal emotion entailment

To explore the causes of emotion expression, early researchers proposed a task called emotion cause extraction (ECE) (Lee, Chen, Li, & Huang, 2010), which aims to extract the reasons behind a certain emotional expression in text. ECE task typically requires emotional expression in advance. Correspondingly, Xia and Ding (2019) proposes the emotion-cause pair extraction task (ECPE) to extract potential emotion-cause pairs in documents and formulate a two-step solution.

**Table 1**

Comparison of existing research methods for CEE and ERC. The commonsense denotes the models inject the commonsense knowledge into network to improve the utterance representation learning.

	Model	(Utterance encoder)	Main method
CEE	RECCON <sup>B;R</sup> (Poria et al., 2021)	Transformer; GCN	Construct the RECCON dataset and use transformer and GCN to capture utterances contextual information.
	TSAM <sup>R</sup> (Zhang et al., 2022)	Attention, GNN	Utilize graph to model speaker and emotion state, then use GAT to interact historical utterances information.
	KEC <sup>R</sup> (Li, Meng et al., 2022)	GNN	Construct a DAG and use knowledge selection strategies to choose commonsense and integrate into utterances.
	KBCIN <sup>R</sup> (Zhao et al., 2023)	GNN	Introduce Commonsense and use the knowledge-enhanced GAT to model the semantic dependencies of the utterances.
	Ours <sup>R</sup> (Jiang, Liu, Tu and Wei, 2023)	Window Transformer	Modeling short-distance relationships between utterances through 2DWindowTransformer.
ERC	DialogueRNN <sup>G</sup> (Majumder et al., 2019)	RNN	Employing RNN to extract dialogic local and global context and modeling speakers.
	DialogueGCN <sup>R</sup> (Ghosal et al., 2019)	GCN	Use GRU to capture context and use GCN to extract speaker-level utterance representation.
	DAG-ERC <sup>R</sup> (Shen et al., 2021)	DAG	Construct a DAG to collect information from neighboring nodes and remote nodes.
	KET <sup>R</sup> (Zhong et al., 2019)	Transformer	Use hierarchical self-attention mechanism and dynamic emotional attention mechanism to model the dialog with integrating commonsense knowledge.
	COSMIC (Ghosal et al., 2020)	RNN	Extract and integrate multi-party commonsense into utterances based on the DialogueRNN model.
	CKCL <sup>R</sup> (Tu, Liang, Mao et al., 2023)	Transformer	The proposed CKCL contrastive learning framework distinguishes which utterances rely on context or external knowledge to identify emotions.
	TODKAT <sup>R;a</sup> (Zhu et al., 2021)	Transformer	Insert a topic layer into a language model, using the Transformer encoder-decoder structure to model topics and fuse commonsense.
	MM-DFN <sup>R</sup> (Hu, Hou, Wei, Jiang, & Mo, 2022)	MGNN	Use GAT to learn context in different semantic spaces, and fuse inter-modal and intra-modal information for conversation.
	DualGATs <sup>R</sup> (Zhang et al., 2023)	DGN	Utilize utterance dependency graphs and speaker dependency graphs to learn contextual information.
	SACL <sup>R</sup> (Hu et al., 2023)	CL	Using the Supervised Adversarial Contrastive Learning framework to learn label-consistent and context-robust emotional features in dialogs.
	Ours <sup>R</sup>	Attention, CEE	Extract utterance emotional cause information through introducing the CEE module and use the reason network in dialog context.

The mark of feature extractors for utterances' context independent representation: B: BERT, G:Glove, R: RoBERTa. Besides, GNN: Graph Neural Network; GCN: Graph Convolutional Network; DAG: Directed Acyclic Graph; MGNN: Multimodal Graph Neural Network; DGN: Dual Graph Network; CL: Contrastive Learning; CEE: Causal Emotion Entailment; The code of this table is available, which can be searched at github.

<sup>a</sup> Denotes the model includes both encoder and decoder.

These tasks all process document data. Unlike emotion cause extraction tasks, the goal of emotion cause entailment is to identify the utterances that trigger the emotion of a specific utterance in a conversation. Previous works, such as Li, Feng, Wang, and Zhang (2019) and Li and Xu (2014), use ECE to solve text-based emotion classification from the perspective of finding emotion-cause and achieved excellent results. Meanwhile, graph construction and transformer (Jiang, Liu, Tu et al., 2023; Li, Li et al., 2022) are used in this task. For example, TSAM (Zhang et al., 2022) uses a graph network to model speakers. Since the dependencies relationships between cause utterances with different emotions from the target utterances are difficult to extract, KEC (Li, Meng et al., 2022) introduces social commonsense knowledge into graph convolution networks to improve the model's reasoning ability for cause utterances. In addition, KBCIN (Zhao et al., 2023) uses commonsense knowledge to build a bridge-interaction network to enhance the understanding of the conversational context. Compared with other conversational tasks, the CEE task is most closely related to conversational emotion recognition at utterance-level. However, no researchers combined CEE with ERC. Table 1 compares mainstream models with our proposed model to highlight their differences.

### 3. Methodology

#### 3.1. Task definition

Let  $U = \{u_1, u_2, \dots, u_N\}$  be a conversation, where  $N$  denotes the utterance quantity. And there is a set  $S = \{S_1, S_2, \dots, S_M\}$  consists of  $M$  speakers. Each utterance  $u_i$  is spoken by the speaker  $S_{\varphi(u_i)}$ , where

$\varphi$  maps the index of the utterance into that of the corresponding its speaker. We also represent  $u_i \in \mathbb{R}^D$  as the feature representation of the utterance. The task of ERC aims to predict the emotion labels of each constituent utterance  $u_i$  from the pre-defined emotion labels (*happy, excited, neutral, angry, sad, frustrated, disgust, fear*). The CEE aims to extract all potential pairs comprised of emotion and corresponding cause clauses from the document annotated with emotion and cause labels in the conversational context. Given a document  $d = [u_1, u_2, \dots, u_i, \dots, u_{|d|}]$ , the purpose of the CEE task is to obtain a series of emotion-cause pairs =  $\{\dots, (u^e, u^{e_1}), \dots, (u^e, u^{e_k}), \dots\}$  where  $u^e$  is an emotion clause and  $u^{e_k}$  is the corresponding  $k$ th cause clause.

#### 3.2. The overall framework

In this section, we present the overall framework: CauERC. The framework consists of the Causal Emotion Entailment (CEE) module, the ERC model, and the Bidirectional Reasoning Network (BRN) module. Six mainstream ERC models are used to test the performance of this framework. Fig. 2 exhibits the whole structure of the presented framework.

#### 3.3. Causal emotion entailment

On account of CEE can correlate emotion-cause with the contextual conversation, we apply the CEE module to the ERC task. Current research focuses more on models based on graph neural networks. For example, TSAM (Zhang et al., 2022) and MPEG (Chen et al., 2023) use

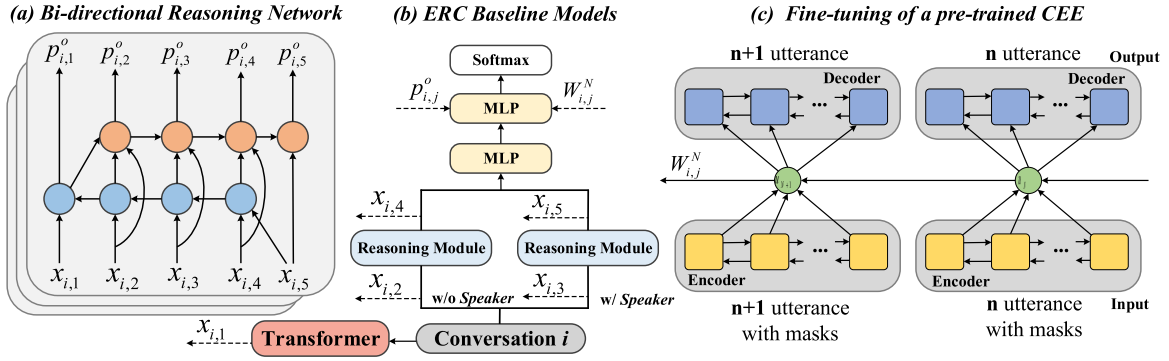


Fig. 2. The overall framework. The part (a) is the BRN module. For a given conversation, we encode the utterances using the Transformer to obtain  $x_{i,1}$  and feed it into the BRN module with the intermediate state vectors  $x_{i,2}, x_{i,3}, x_{i,4}, x_{i,5}$  obtained from the ERC model. Part (c) is a pre-trained model. Part (b) is the ERC model, and the part of the ERC model that extracts the utterance context representation is used as a whole as the contextual reasoning module. We concatenate  $x_{i,4}$  with  $x_{i,5}$  and feed it into the first MLP. The dotted lines indicate the direction of data propagation from different modules to each other. The solid lines mean the direction of information propagation between different nodes within each module, whereas the solid lines in part (c) indicate the process of encoding and decoding the words in the utterance.

attention mechanisms and graph networks to fuse speaker and sentiment information. KEC (Li, Meng et al., 2022) and KBCIN (Zhao et al., 2023) introduce commonsense knowledge and use graph neural networks. However, they perform poorly in information exchange between sentences at short distances. Inspired by our previous work (Jiang, Liu, Tu et al., 2023), the emotional cause information is sensitive to the position of clauses in conversation. We use window transformers to extract contextual information within short distance when handling CEE tasks. Rather than modeling the entire conversation or injecting external information to enhance utterance representations, we focus on the impact of inter-utterance interaction in a short distance, and use 2D window transformer to interact inter-utterances semantic information with limited window size. Specifically, the 2D Window Transformer is used in the pre-trained process. The given utterances are divided into several windows according to window size to get better clause representation with modeling the relationship between clauses. This module has  $N$  encoder layers. Each layer comprises a window attention and a feed-forward layer. The 2D Window Transformer is utilized as the encoder layer of the CEE module. Each utterance pair  $(u_i, u_j)$  is fed into embedding layer to get the representation  $W_{i,j}$ . Firstly,  $W_{i,j}$  is calculated by window attention which is multi-head self-attention. The  $W_{i,j}$  is fed into three linear layers to calculate the query vector  $q_{i,j}$ , key vector  $k_{i,j}$  and the value vector  $v_{i,j}$ .

$$q_{i,j} = W_{i,j} W_q \quad (1)$$

$$k_{i,j} = W_{i,j} W_k \quad (2)$$

$$v_{i,j} = W_{i,j} W_v \quad (3)$$

where  $W_q \in \mathbb{R}^{n \times n}$ ,  $W_k \in \mathbb{R}^{n \times n}$  and  $W_v \in \mathbb{R}^{n \times n}$  are learned parameters. For the three vectors  $q_{i,j}$ ,  $k_{i,j}$  and  $v_{i,j}$ , the weight  $\beta_{i,j}$  and the output of window attention is calculated as follows:

$$\beta_{i,j} = \text{softmax} \left( \frac{k_{i,j}^T \cdot q_{i,j}}{\sqrt{n}} \right) \quad (4)$$

$$z_{i,j} = v_{i,j} \beta_{i,j}^T \quad (5)$$

where  $z_{i,j}$  is the output of window attention. The input for the feed-forward layer is  $z_{i,j}$  entered into a layer that has two identical constructions followed by a normalization layer at its output:

$$o_{i,j,1} = \text{dropout} (z_{i,j} W_1 + b_1) \quad (6)$$

$$o_{i,j,2} = o_{i,j,1} + \text{dropout} (o_{i,j,1} W_2 + b_2) \quad (7)$$

$$o_{i,j} = o_{i,j,2} + \text{norm} (o_{i,j,2}) \quad (8)$$

where the *norm* denotes laynorm layer.  $o_{i,j,1}$  and  $o_{i,j,2}$  are the output of the two sublayers, respectively.  $o_{i,j}$  is the output of a encoder layer in 2D Window Transformer.

$$W_{i,j}^{t+1} = o_{i,j}^t \quad (9)$$

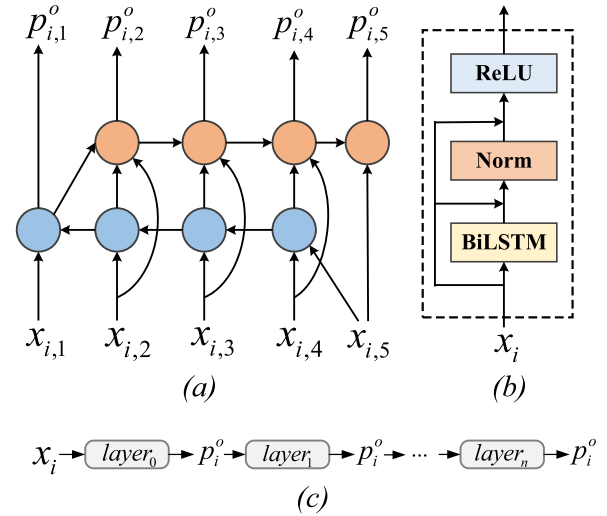


Fig. 3. The Bidirectional Reasoning Network (BRN) module, which have several layers. (a) is the overall structure of sublayer, and (b) is the structure of each node in BRN layer. The subfigure (c) show the iterative process between layers. The  $x_i$  and  $p_i$  are inputs and outputs. The input  $x_{i,n}$ , ( $n \in [1, \dots, 5]$ ) represents the utterances representation of the model at different stages. In subfigure c, the  $p_i^o$  is the middle representation inter-sublayers, which is the output of last sublayer while the input of next sublayer.

where the output  $W_{i,j}^N$  of the last layer is the representation of utterance pair  $(u_i, u_j)$  extracted by 2D Window Transformer. The relative position modeling is used to learn the representation of clauses pair and ranks the candidate clauses.

By saving the pre-trained weight and transferring this model, we can convert low-level clause representation to high-level representation, which contains information about the cause evoking the clause. The utterance representation extracted by the pre-trained model is expressed as  $md_i$ .

#### 3.4. Bidirectional reasoning network

In a conversation, there is emotional dependence and emotional influence of the speaker. People emotional cognition and reasoning processes will affect emotional dependence (Lerner et al., 2015). To simulate human emotional and cognitive behavior in dynamic conversations, emotion recognition models need to perform complex reasoning and learning with context. In this section, we design the BRN module to capture the context of emotional information.



Instead of employing BiLSTM to extract contextual information, each node in the BRN structure contains a BiLSTM. Each node extracts the semantic information of the utterances and fuses the utterances representations of different stages of the model through a bidirectional pathway structure. Through multi-turn iteration, BRN can simulate the human emotional reasoning process. The BRN module has  $N$  layers and the structure of a layer is shown in Fig. 3(a). Specially, the input of BRN module is  $(x_{i,1}, x_{i,2}, x_{i,3}, x_{i,4}, x_{i,5})$ , and there are two pathways for information fusion. In the right-to-left pathway, the output  $p_{i,k}^{rl}$  of each node is calculated as follows:

$$p_{i,k}^{rl} = \text{cell} \left( w_{k,1} \cdot x_{i,k} + w_{k,2} \cdot p_{i,k+1}^{rl} \right) \quad (10)$$

where  $p_{i,5}^{rl} = x_{i,5}$  and  $w_{k,i}$  is trainable weight that can be a scalar.  $\text{cell}$  is the node of the BRN module. In the left-to-right pathway, the output of each node is calculated as follows:

$$p_{i,k}^o = \text{cell} \left( w'_{k,1} \cdot x_{i,k} + w'_{k,2} \cdot p_{i,k}^{rl} + w'_{k,3} \cdot p_{i,k-1}^o \right) \quad (11)$$

where  $p_{i,1}^o = x_{i,1}$  and  $w'_{k,i}$  is trainable weight that can be a scalar.  $k$  is the index of the cell in the BRN module. The Fig. 3(b) exhibits the structure of  $\text{cell}$ . This utterance representation is fed into BiLSTM, which is followed by the norm layer and activation layer. The output of these cells for the input  $x_i$  can be computed as:

$$c_i = \text{norm}(x_i + \text{BiLSTM}(x_i)) \quad (12)$$

$$p_i^o = \text{ReLU}(x_i + c_i) \quad (13)$$

where the norm layer is LayerNorm and the activation layer uses the ReLU function. The output of the current layer is the input of the next layer.

$$x_{i,k}^{l+1} = p_{i,k}^l \quad (14)$$

where  $l$  is the index of the layer. Then, the output of the whole reasoning network is obtained by concatenating the final layer's output. In general, given  $(x_{i,1}, x_{i,2}, x_{i,3}, x_{i,4}, x_{i,5})$ , the vector  $mf_i$  extracted by of BRN module can be defined as:

$$mf_i = \text{BRN}(x_{1,i}, x_{2,i}, x_{3,i}, x_{4,i}, x_{5,i}) \quad (15)$$

where the  $mf_i$  is the output of the BRN module.

### 3.5. Emotion classifier

Based on the output vectors  $md_i, mf_i$  obtained from the BRN module and the CEE module, respectively. We concatenate them with the vector  $mc_i$  obtained from the ERC model's last layer and fuse them using MLP to gain the utterance representation  $o_i$ .

$$o_i = \text{MLP}([md_i; mf_i; mc_i]) \quad (16)$$

where the  $o_i$  is the final representation fed into the emotion classification layer employed for emotion prediction:

$$\hat{y}_i = \text{softmax}(W_o o_i + b_o) \quad (17)$$

The cross-entropy loss function is applied to calculate the loss value to optimize the model:

$$\text{loss} = - \frac{1}{\sum_{l=1}^L c(l)} \sum_{i=1}^L \sum_{k=1}^{c(i)} y_{i,k}^l \log(\hat{y}_{i,k}^l) \quad (18)$$

where  $L$  is the number of the conversation.  $c(l)$  denotes the number of utterance in the conversation  $i$ .  $y_{i,k}^l$  and  $\hat{y}_{i,k}^l$  are the true label of utterance  $i$  in conversation  $l$  and the possibility of predicting the result of category  $k$ , respectively.

**Table 2**

The statistics of datasets. Statistics of splits and evaluation metrics used in different datasets.

Datasets	Conversations			Utterances		
	Train	Val	Test	Train	Val	Test
IEMOCAP	120	12	31		5810	1623
DailyDialog	11,118	1000	1000	87,832	7912	7863
Datasets	Classes			Evaluation		
IEMOCAP	6			Accuracy and Weighted F1		
DailyDialog	7			Macro F1 and Micro F1		

## 4. Experimental settings

### 4.1. Datasets and evaluation metrics

Our framework is evaluated on the following datasets: IEMOCAP (Busso et al., 2008), DailyDialog (Li et al., 2017). The detailed statistics of the datasets are reported in Table 2.

- **IEMOCAP** is a multimodal dataset for emotion recognition that is comprised of videos of multi-turn dialogs of ten unique speakers. The utterances are annotated with one of six emotion labels, namely *happy, excited, neutral, angry, sad, and frustrated*.
- **DailyDialog** is an emotion detection dataset that contains the conversations of our daily life and human-written daily communications. There are seven emotion labels annotated by three professional persons in this dataset: *disgust, fear, sadness, angry, neutral, joy, surprise*.

Because of the uneven distribution of the DailyDialog dataset, the percentage of utterances with the neutral label is 83%, so we adopt the Micro F1 and Macro F1, excluding the neutral samples. We follow the previous research (Majumder et al., 2019) to use average Accuracy (Acc.) and Weighted F1 on the IEMOCAP dataset. In this paper, we leverage MaF to represent the Macro F1, WF to represent the Weighted F1, and MiF to represent the Micro F1, respectively.

### 4.2. Baselines

To evaluate the performance of the proposed framework, we compare it with several baselines.

1. **KET** (Zhong et al., 2019) introduces external common-sense knowledge into a transformer architecture through self-attention and graph-attention mechanisms.
2. **VHRED** (Hazarika, Poria, Zimmermann, & Mihalcea, 2021) uses a pre-trained sentence encoder and simulates the inter-sentence context through transfer learning to identify the emotion.
3. **DialogueRNN** (Majumder et al., 2019) models context and speaker separately using GRU to obtain global context dependencies and speaker dependencies, meanwhile using global GRU for speaker-to-speaker interaction.
4. **DialogueGCN** (Ghosal et al., 2019) models different speakers using GCN pairs after capturing contextual information separately and classifies the emotions of utterance representations by attention mechanism.
5. **BiERU** (Li, Shao, Ji and Cambria, 2022): construct a bidirectional sentiment recursive unit by utilizing many GRU to detect emotion.
6. **RoBERTa** (Zhang, Xu, & Dong, 2020) uses the pre-trained RoBERTa to obtain utterance representation and fine-tune the prediction layer.
7. **COSMIC** (Ghosal et al., 2020) employs RoBERTa (Zhang et al., 2020) to extract the data of this paper and introduce commonsense knowledge like mental state, causality, etc. Using the pre-trained model COMET (Bosselut et al., 2019) and feed them into the emotion analysis model.

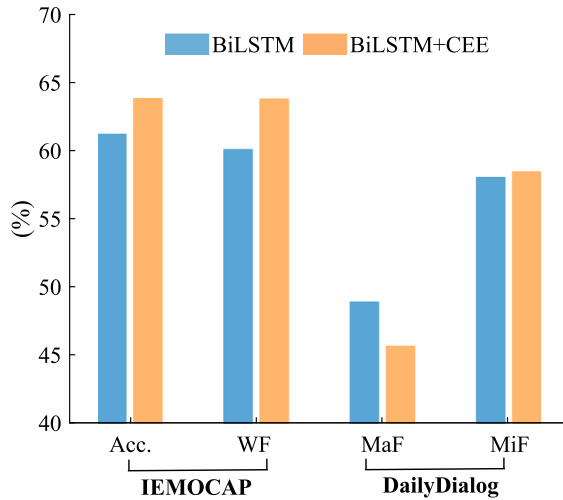


Fig. 4. Experimental results for the verification of CEE validity on two different datasets.

8. **DialogueCRN** (Hu, Wei, & Huai, 2021) processes utterances representation by using BiLSTM and attention mechanism to simulate the human cognitive.
9. **SKAIG** (Li, Lin, Fu, & Wang, 2021) introduces commonsense knowledge into the graph structure.
10. **DAG-ERC** (Shen et al., 2021) combines traditional graph-based models with recursive-based neural models.
11. **MM-DFN** (Hu et al., 2022) leverages graph-base construction learning the intra- and inter-modal relationship of utterances.

In this paper, we employ the pre-trained model 840B GloVe (Pennington, Socher, & Manning, 2014) to obtain the utterance representation with a dimension of 300. The extracted utterance representations are then fed into a network consisting of a convolutional layer, maximum pooling, and fully connected layers to extract text features. The final vector with a dimension of 100 is used as the text feature.

In addition to using GloVe (Pennington et al., 2014) as the feature extractor, the pre-trained model BERT (Kenton & Toutanova, 2019) and RoBERTa (Zhang et al., 2020) is also applied to extract context-independent text features. The output vector of the final layer of the pre-trained model is used as the text feature.

#### 4.3. Hyperparameters settings

We conduct hyperparameters search for our proposed framework on IEMOCAP and DailyDialog datasets. We employ Adam optimization with a batch size of 32, epochs of 50, the learning rate of  $\{1e-5, 2e-5\}$ , L2 weight decay of  $2e-4$ , and dropout of  $\{0.3, 0.2\}$ . The number of 2D Window Transformer's encoder layers is 3 and the window size is 4 in the CEE module. The number of layers in the BRN module is 2.

## 5. Results and discussions

### 5.1. The role of causal emotion

To prove the validity of CEE in the ERC model, we conduct one of the most classical models BiLSTM for emotion recognition, where the RoBERTa is applied to extract the textual features. The results are displayed in Fig. 4. Compared with using BiLSTM only for emotion classification, the experimental results of BiLSTM+CEE are improved. This improvement underscores how incorporating CEE enhances the model's ability to capture causal emotional cues, thereby strengthening its overall effectiveness.

Table 3

The experimental results. The results in bold are the best-performing ones under each column. The best values are highlighted in bold. All the results of the comparable baselines can be found in papers (Hu et al., 2022; Shen et al., 2021).

Models	IEMOCAP		DailyDialog		
	Acc.	WF	MaF	MiF	
GloVe-based	1. KET	-	59.56	-	53.37
	2. VHRED	-	58.60	-	48.40
	3. DialogueRNN	63.03	62.50	-	50.56
	4. DialogueGCN	65.25	64.18	-	-
	5. BiERU	63.02	63.13	-	-
	6. DialogueCRN	65.25	65.21	-	-
RoBERTa-based	7. RoBERTa	-	63.38	48.20	55.16
	8. COSMIC	-	65.28	51.05	58.48
	9. SKAIG	-	66.98	51.95	59.75
	10. DAG-ERC	-	68.08	-	59.33
	11. MM-DFN	68.21	68.18	-	-
	<b>CauERC</b>	<b>69.01</b>	<b>69.07</b>	<b>53.39</b>	<b>59.81</b>

Table 4

The experimental results of generalization analysis. The best values are highlighted in bold.

Models	IEMOCAP		DailyDialog	
	Acc.	WF	MaF	MiF
BiERU	63.22	63.52	29.35	52.79
<b>CauERC+BiERU</b>	62.57	62.45	<b>39.30</b>	<b>56.37</b>
DialogueRNN	64.20	64.21	39.69	56.19
<b>CauERC+DialogueRNN</b>	<b>66.42</b>	<b>66.37</b>	<b>51.29</b>	<b>58.59</b>
DialogueCRN	66.54	66.11	52.25	58.28
<b>CauERC+DialogueCRN</b>	<b>69.01</b>	<b>69.07</b>	<b>53.39</b>	<b>59.81</b>

Table 5

The results of significance tests for generalization analysis (P-Value).

Datasets	BiLSTM	BiERU	DialogueRNN	DialogueCRN
IEMOCAP	9.74e-3	7.84e-3	9.46e-3	3.67e-3
DailyDialog	4.57e-6	6.58e-6	9.85e-3	1.67e-2

P-Values < 0.05.

### 5.2. Experimental results and analysis

Our framework based on the DialogueCRN model is compared with the baselines in Table 3. To obtain our results in Table 3, we employ RoBERTa to extract text features. As expected, our framework outperforms all the baselines. On the IEMOCAP dataset, we achieve a new state-of-the-art Acc. of 69.01% and WF of 69.07%. Compared with the previous work, our framework gains 0.80%, and 0.89% in terms of Acc. and WF. On the DailyDialog dataset, our framework gets a 1.44% and 0.06% improvement on MaF and MiF.

In order to explain the gaps in experimental results, it is essential to understand the logical relationships of the conversations. Previous works focus on model speakers or introducing external knowledge to enrich contextual representations. They both encode the utterances, but none of them consider the connection between utterance and its associated cause clauses. In the process of contextual information propagation, the ERC model gradually loses the information between the cause clauses associated with the current utterance. In contrast, we solve these two problems with the CEE and BRN modules.

### 5.3. Generalization analysis

Based on the promising results of the CEE+BiLSTM model, we extended our framework to the DialogueCRN model to evaluate the effectiveness of combining the CEE module and BRN module with the ERC model. To assess the generalizability of our framework, we compared the results of two emotion recognition models, BiERU (Li, Shao et al., 2022) and DialogueRNN (Majumder et al., 2019), as shown in Table 4.

**Table 6**  
The experimental results of ablation studies on IEMOCAP and Dailydialog datasets.

	BRN	CEE	IEMOCAP						DailyDialog					
			$\mathfrak{R}_{BERT}$			$\mathfrak{R}_{RoBERTa}$			$\mathfrak{R}_{BERT}$			$\mathfrak{R}_{RoBERTa}$		
			Acc.	MaF	WF	Acc.	MaF	WF	Acc.	MaF	MiF	Acc.	MaF	MiF
RoBERTa-based	×	×	66.54	65.81	66.11	66.54	65.81	66.11	83.13	<b>52.35</b>	58.28	83.13	52.35	58.28
	×	✓	67.71	66.81	67.43	68.33	66.07	67.97	84.33	50.32	58.37	85.12	49.01	58.38
	✓	×	67.90	<b>67.11</b>	68.19	67.90	67.11	68.19	80.13	49.78	56.27	80.13	49.78	56.27
	✓	✓	<b>68.86</b>	67.04	<b>68.58</b>	<b>69.01</b>	<b>68.15</b>	<b>69.07</b>	<b>85.63</b>	51.10	<b>59.30</b>	<b>85.56</b>	<b>53.39</b>	<b>59.81</b>

$\mathfrak{R}$  represents the CEE Pre-trained Model. The features used in the ERC task are extracted by RoBERTa. The best values are highlighted in bold.

The results of Table 3 indicate that RoBERTa outperforms GloVe in extracting textual features. Therefore, in this section, we use RoBERTa to extract text features. Overall, the experimental results demonstrate the effectiveness of our framework. The results of the significance tests on the compared models are reported in Table 5, which demonstrates that our framework is significantly different from the comparison models. Moreover, in order to elucidate the usage of RoBERTa as a feature extractor and explore the adaptability of the proposed framework to ERC tasks, we employed three pre-trained models, namely GloVe, BERT, and RoBERTa, to extract text representations. Extensive comparative experiments were conducted on frameworks utilizing these three extractors, with the DialogueCRN model employed in the experiments. The experimental results, as depicted in Fig. 5, indicate that the RoBERTa-based framework outperforms the others, suggesting that the representations derived from a more powerful extractor, such as RoBERTa, yield greater benefits for emotional recognition.

#### 5.4. Ablation study

To investigate the contribution of the proposed modules, we conducted several ablation studies on the DialogueCRN model, where each constituent component was removed individually. As shown in Table 6, the performance decreases slightly when either the BRN or CEE module is removed, indicating the significance of both modules.

**Analysis of Bidirectional Reasoning Network:** As demonstrated in Table 6, when RoBERTa is employed to extract text features, the results are enhanced on the IEMOCAP dataset, with the evaluation indicators Acc. and WF increasing by 1.36% and 2.08%, respectively. Unlike the reasoning network in DialogueCRN, which directly processes utterance representations through lstm-attention, while Bidirectional Reasoning Network (BRN) has a bidirectional structure. BRN does not simply act as a neural network at a certain layer of the model but dynamically processes different levels of utterance representation of the emotion recognition model to conduct emotional reasoning. These results showcase the effectiveness of the bidirectional reasoning network in extracting emotion cues by utilizing the representation from the intermediate layers of the ERC model. Moreover, it successfully simulates the process of human-like emotional reasoning in conversations through multiple iterations, enhancing the overall comprehension and coherence of the framework.

**Analysis of Causal Emotion Entailment:** Table 6 presents the results of our experiments, which show that using CEE in the DialogueCRN model improves the performance, regardless of whether BERT or RoBERTa is used as the text feature extractor. Specifically, when RoBERTa is used, the evaluation metrics Acc. and WF on the IEMOCAP dataset show an improvement of 1.79% and 1.86%, respectively, with values of 68.33% and 67.97%. However, the results of MaF and MiF on the DailyDialog dataset show weak improvement, with values of 49.01% and 58.38%, respectively. This indicates that focusing on the corresponding cause information through the CEE module can enhance the emotion reasoning ability of the conversation emotion recognition model, particularly in multi-round conversations. In the case of combining BRN with CEE, the result is better than that using one of them alone on IEMOCAP and DailyDialog datasets.

We can conclude that ablating both modules simultaneously results in a greater decline in ERC model performance. This indicates that the BRN module and the CEE module can complement each other. Although our method is evaluated on the IEMOCAP and DailyDialog datasets, it is not constrained to conversations with two participants. According to our ablation analysis, the BRN proposed in this paper utilizes the model's multiple intermediate layer outputs to extract emotional cues and simulate human emotional cognitive behavior. The CEE module enables the model to focus on relevant causal information, enhancing its performance. Importantly, both modules do not limit the number of speakers in a conversation. Table 4 also presents the results of applying our method to several emotion detection models, which can be used in dialog systems involving multiple speakers. To better adapt to multi-speaker conversations, it is essential to annotate corresponding labels in the multi-speaker conversation dataset. This may pose some challenges, such as emotional dependency within conversations, which is the emotional interactions between speakers. We can consider fusing personalized factors into sentence representation extraction.

#### 5.5. Case study

In this section, we present a case study on a conversation example from the Dailydialog dataset in Fig. 6, which shows the cause clause's role in the utterance's emotion. To validate the efficacy of our framework, two individuals were tasked with annotating emotion-cause clause labels for each utterance in the case study samples, indicating whether the utterance contains causal information. The connecting lines in the graph depict the emotion-cause relationships between the utterances. Furthermore, we visualized the attention layer of the final ERC model for enhanced comprehension. The emotion label of **utterance 5** and **utterance 7** is easily predicted to be *neutral*, while the actual label is *happiness*. As shown by the solid lines in Fig. 6, the cause clauses of **utterance 5** are **utterance 1** and **utterance 3**. At the same time, the dotted lines indicate that the cause clauses of **utterance 7** are **utterance 3**, **utterance 5**, and **utterance 6**.

When not performing the CEE, **utterance 5** and **utterance 7** obtain contextual information by using the ERC model, we will incorrectly predict their emotion as *neutral* after a few rounds of training. In contrast, with the CEE task, the model enhances the effect of the cause clause on the associated utterance emotion. During model training, the contextual representation of **utterance 5** will contain more information related to **utterance 1** and **utterance 3**, and the contextual representation of **utterance 7** will contain more information related to **utterance 3**, **utterance 5**, and **utterance 6**, which causes **utterance 5** and **utterance 7** to be correctly predicted as *happiness*. The illustrated utterances show that CEE highlights the influence of some clauses on a particular utterance.

#### 5.6. Emotion interpretability analysis

In this section, we conduct an emotion interpretability analysis on a conversation and display the visualization of each utterance. The conversation is sampled from the conversational dataset Dailydialogue. We employ the model interpretability tool LIME (Ribeiro, Singh, & Guestrin, 2016) to conduct the emotion interpretability analysis.

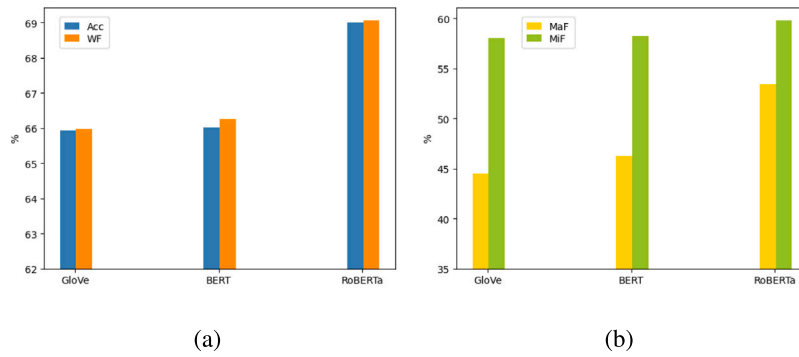
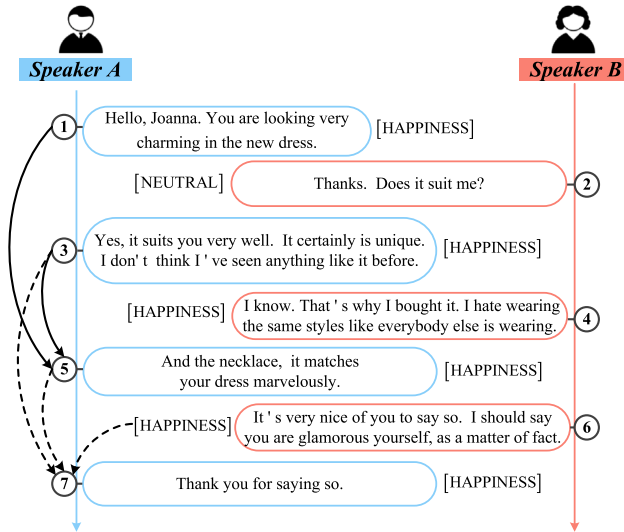
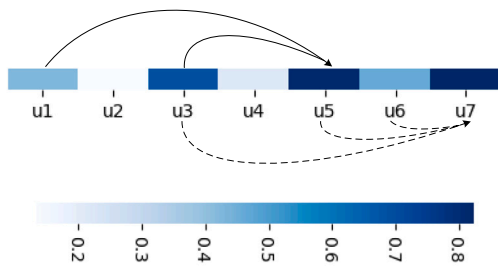


Fig. 5. Performance on different feature extractors. (a) IEMOCAP. (b) DailyDialog.



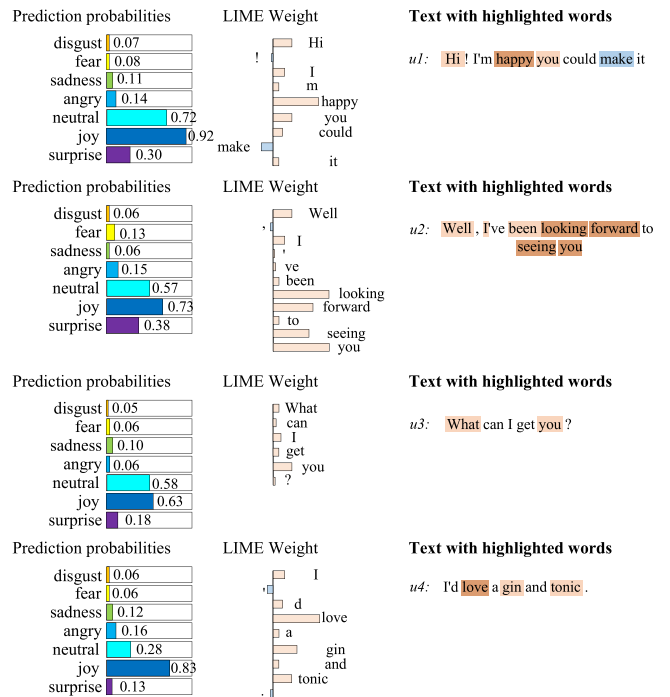
(a) A conversation from DailyDialog dataset.



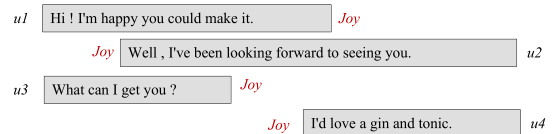
(b) The visualization of attention weights of the utterances in conversation. Each cell represents a utterance ( $u_i$ ).

Fig. 6. Case study of a conversation from the DailyDialog.

Fig. 7 demonstrates the analysis result. The four utterances in this conversation are annotated with the emotion label of joy. The prediction probabilities for the neutral and joy categories are higher, which may be due to the large proportion of neutral samples in the DailyDialog dataset. Additionally, emotions opposite to joy, such as disgust and fear, have lower prediction probabilities, which further suggests the influence of emotional inertia in conversations on emotion analysis. By examining the weights of the highlighted words in the text, we can observe that the words influencing the emotion of the utterances are typically among the more important ones. These words entail the cause information for the emotion clauses.



The relationship between cause clauses and emotion clauses in this conversation.



External Knowledge  
 xIntent: want to {have a drink / quench thirst / drink something sweet}.  
 xNeed: need to {go to the bar / buy gin and tonic / get a drink}.  
 xReact: feels {happy / satisfied / relaxed}.  
 xWant: want to {have a drink / drink the gin / order a drink}.  
 xEffect: to {gets drunk / drinks a lot / drinks the gin}  
 oReact: feels {happy / thirsty / excited}.

Fig. 7. The emotion interpretability analysis of a short conversation in DailyDialog dataset.

For example, in utterances  $u_1$ ,  $u_3$ , and  $u_4$ , the utterances  $u_1$  and  $u_3$  plays a guiding role for  $u_4$ , thereby affecting  $u_4$  emotion. Meanwhile, the highlighted words and their corresponding weights reveal a strong correlation between utterances. Each utterance in the dialog affects each other.



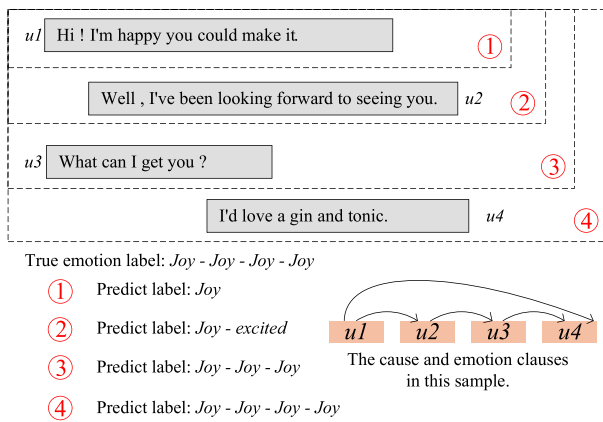


Fig. 8. The emotion recognition analysis in a dynamic conversation on DailyDialog dataset. There are four test results of the sample in the figure.

Additionally, we have also extracted the commonsense knowledge for the utterances of sample through the commonsense model comet, which can reflect the human action and psychological activities. In Fig. 7, when the person say he love gin and tonic, the extracted commonsense display he wants to have a drink and he is happy. The result can be inferred from the context and the causal clauses  $u1$  and  $u3$ . The sample visualization result indicates that the model can simulate human emotional cognition by reasoning on the contextual semantic information extracted from the conversation.

### 5.7. Dynamic emotion recognition study

When constructing a emotion analysis model in conversational system, it is crucial to analyze emotional changes in dynamic conversations in real-time. The CEE task is used to extract the entailed cause information contained in the clause, which utilizes the utterance information produced before the current utterance. Additionally, when identifying the emotion of an utterance, we can input this utterance into the model together with the preceding utterances, and then use the proposed framework to reason and identify the emotion labels through these utterances. In this section, we have explored the influence of the temporal dimension of dynamic conversations on emotion analysis. The dialog samples are fed into the model for sentiment label prediction. The visualization results of dynamic emotion prediction are shown in Fig. 8.

When two utterances are input into the model, the utterance  $u2$  predicted emotion label is *excited*, which indicates that when the number of input utterances is limited, the model lacks sufficient context and causal information for emotion inference. As the number of input utterances increases, the available information to the model gradually increases, which improves the model's emotion recognition performance. The result suggests that the model can capture conversational temporal dependencies by leveraging context. Even so, there is still room for improvement in the model's emotion recognition performance in dynamic conversation. Fig. 8 shows a short conversation with four utterances. In long conversations, rich context can provide more information, and enhancing the scalability of the model requires a balance between computational efficiency and context capture capabilities. We can explore dynamic context and cause information selection mechanisms based on reinforcement learning to filter important historical fragments, avoid interference from irrelevant information, and further optimize resource allocation. Meanwhile, the performance of model will also be improved with commonsense enhanced. In addition, when considering the application of the proposed framework in practice, the framework needs to be adaptable to different scenarios.

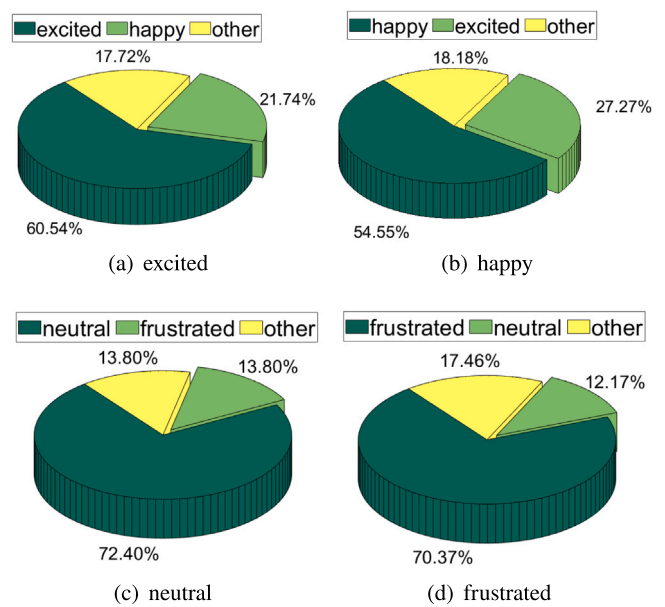


Fig. 9. Experimental results of error analysis. Four similar emotion prediction results on the IEMOCAP dataset. For example, excited is misclassified as happy and other. Its percentages are as shown in the subgraph (a) above.

This method is evaluated on the datasets IEMOCAP and Dailydialogue, which are English dialog datasets. When adapting it to new language scenarios, we can fine-tune on cross-lingual emotion recognition datasets in future work, apply large language model to fine-tune or use some domain adaptation techniques. Concurrently, constructing the causal graphs through extracting the emotion causal information and leveraging the counterfactual reasoning method to improve conversational causal reasoning and interpretability, so that the framework can be better used in mental health and emotional support or conversational agents. When considering the model deployment, we also need to edge computing and federated learning to meet the real-time requirements and data privacy of the application.

### 5.8. Error analysis

Although our framework has shown strong performance, it still has some limitations. The analysis of our experimental results indicates that our model struggles to effectively distinguish between similar emotion categories such as excited, happy, neutral, and frustrated. Fig. 9 illustrates the classification results of our experiments on these four emotion categories in the IEMOCAP dataset. A similar situation is observed in the DailyDialog dataset. We suspect that this difficulty arises because utterances with similar emotions have similar semantic information in the extracted features.

Furthermore, our experiments are limited to text data, whereas multimodal data can provide additional information for non-neutral emotions in utterances. For example, videos may show a disappointed expression for utterances with sad emotions, and utterances with angry emotions may have a higher pitch. However, due to the limitations of the CEE task, not all utterances used for emotion analysis can obtain information about their corresponding cause clauses, which ultimately limits the performance of our framework.

## 6. Conclusions and future work

This paper proposes a framework that combines CEE and BRN to enhance the ability of emotion analysis in ERC. Specifically, our framework emphasizes the causal clause that triggers emotions via the CEE module and addresses the issue of disregarding the context of

other clauses when the CEE module is integrated with the ERC model, with the help of the BRN module. The proposed framework achieves state-of-the-art result on two public conversational emotion recognition datasets. Nevertheless, there are still some shortcomings, such as the limitation of datasets in CEE tasks. Therefore, we plan to annotate the dataset by taking context and personal relationships into account in future work. And apply the proposed method to more conversational scenarios. Besides, we will pay attention to incorporating multimodal information into this framework and effectively fusing commonsense information.

### CRedit authorship contribution statement

**Hao Liu:** Conceptualization, Data curation, Writing – original draft, Software. **Runguo Wei:** Validation. **Geng Tu:** Investigation. **Jiali Lin:** Methodology. **Dazhi Jiang:** Supervision. **Erik Cambria:** Writing – review & editing.

### Ethical approval

This article does not contain any studies with human participants or animals performed by any of the authors.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgments

This research/project is supported by the Ministry of Education, Singapore under its MOE Academic Research Fund Tier 2 (STEM RIE2025 Award MOE-T2EP20123-0005) and by the National Natural Science Foundation of China (62372283, 62206163), Science and Technology Major Project of Guangdong Province, China (STKJ2021005, STKJ202209002, STKJ2023076), Natural Science Foundation of Guangdong Province (2024A1515010239), and Key Project of Colleges and Universities in Guangdong Province (Natural Science), China (2024ZDZX1019).

### Data availability

No data was used for the research described in the article.

### References

- Bosselut, A., Rashkin, H., Sap, M., Malaviya, C., Celiyilmaz, A., & Choi, Y. (2019). Comet: Commonsense transformers for automatic knowledge graph construction. In *Proceedings of the 57th annual meeting of the association for computational linguistics* (pp. 4762–4779).
- Busso, C., Bulut, M., Lee, C.-C., Kazemzadeh, A., Mower, E., Kim, S., et al. (2008). IEMOCAP: Interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42(4), 335–359.
- Cambria, E., Howard, N., Hsu, J., & Hussain, A. (2013). Sentic blending: scalable multimodal fusion for the continuous interpretation of semantics and sentics. In *Proceedings of IEEE SSCI* (pp. 108–117).
- Cambria, E., Mao, R., Chen, M., Wang, Z., & Ho, S.-B. (2023). Seven pillars for the future of artificial intelligence. *IEEE Intelligent Systems*, 38(6), 62–69.
- Cambria, E., Zhang, X., Mao, R., Chen, M., & Kwok, K. (2024). SenticNet 8: Fusing emotion AI and commonsense AI for interpretable, trustworthy, and explainable affective computing. In *Proceedings of the International Conference on Human-Computer Interaction* (pp. 197–216). Washington DC, USA.
- Chen, T., Shen, Y., Chen, X., Zhang, L., & Zhao, S. (2023). MPEG: A multi-perspective enhanced graph attention network for causal emotion entailment in conversations. *IEEE Transactions on Affective Computing*, 1–14.
- Chudasama, V., Kar, P., Gudmalwar, A., Shah, N., Wasnik, P., & Onoe, N. (2022). M2FNet: Multi-modal fusion network for emotion recognition in conversation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 4652–4661).
- Ding, Z., Xia, R., & Yu, J. (2020). ECPE-2D: emotion-cause pair extraction based on joint two-dimensional representation, interaction and prediction. In *Proceedings of the 58th annual meeting of the association for computational linguistics* (pp. 3161–3170).
- Fan, C., Lin, J., Mao, R., & Cambria, E. (2024). Fusing pairwise modalities for emotion recognition in conversations. *Information Fusion*, 106, Article 102306.
- Ghosal, D., Majumder, N., Gelbukh, A., Mihalcea, R., & Poria, S. (2020). Cosmic: Commonsense knowledge for emotion identification in conversations. In *Findings of the association for computational linguistics: EMNLP 2020* (pp. 2470–2481).
- Ghosal, D., Majumder, N., Poria, S., Chhaya, N., & Gelbukh, A. (2019). DialogueGCN: A graph convolutional neural network for emotion recognition in conversation. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing* (pp. 154–164).
- Hazarika, D., Poria, S., Gorantla, S., Cambria, E., Zimmermann, R., & Mihalcea, R. (2018). CASCADE: Contextual sarcasm detection in online discussion forums. In *Proceedings of the 27th international conference on computational linguistics* (pp. 1837–1848).
- Hazarika, D., Poria, S., Mihalcea, R., Cambria, E., & Zimmermann, R. (2018). ICON: Interactive conversational memory network for multimodal emotion detection. In *Proceedings of the 2018 conference on empirical methods in natural language processing* (pp. 2594–2604).
- Hazarika, D., Poria, S., Zimmermann, R., & Mihalcea, R. (2021). Conversational transfer learning for emotion recognition. *Information Fusion*, 65, 1–12.
- Hu, D., Bao, Y., Wei, L., Zhou, W., & Hu, S. (2023). Supervised adversarial contrastive learning for emotion recognition in conversations. In A. Rogers, J. Boyd-Graber, & N. Okazaki (Eds.), *Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers)* (pp. 10835–10852). Toronto, Canada: Association for Computational Linguistics, <http://dx.doi.org/10.18653/v1/2023.acl-long.606>, URL: <https://aclanthology.org/2023.acl-long.606>.
- Hu, D., Hou, X., Wei, L., Jiang, L., & Mo, Y. (2022). Mm-dfn: Multimodal dynamic fusion network for emotion recognition in conversations. In *ICASSP 2022-2022 IEEE international conference on acoustics, speech and signal processing* (pp. 7037–7041). IEEE.
- Hu, D., Wei, L., & Huai, X. (2021). DialogueCRN: Contextual reasoning networks for emotion recognition in conversations. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing (volume 1: long papers)* (pp. 7042–7052).
- Huang, S., Liu, Y., Tsang, I. W., Xu, Z., & Lv, J. (2022). Multi-view subspace clustering by joint measuring of consistency and diversity. *IEEE Transactions on Knowledge and Data Engineering*.
- Jia, B., Lv, J., Peng, X., Chen, Y., & Yang, S. (2020). Hierarchical regulated iterative network for joint task of music detection and music relative loudness estimation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29, 1–13.
- Jiang, D., Liu, H., Tu, G., & Wei, R. (2023). Window transformer for dialogue document: a joint framework for causal emotion entailment. *International Journal of Machine Learning and Cybernetics*, 1–11.
- Jiang, D., Liu, H., Tu, G., Wei, R., & Cambria, E. (2024). Self-supervised utterance order prediction for emotion recognition in conversations. *Neurocomputing*, 577, Article 127370.
- Jiang, D., Liu, H., Wei, R., & Tu, G. (2023). Csat-ftcn: A fuzzy-oriented model with contextual self-attention network for multimodal emotion recognition. *Cognitive Computation*, 1–10.
- Jiang, D., Wei, R., Liu, H., Wen, J., Tu, G., Zheng, L., et al. (2021). A multitask learning framework for multimodal sentiment analysis. In *2021 international conference on data mining workshops* (pp. 151–157). <http://dx.doi.org/10.1109/ICDMW53433.2021.00025>.
- Kenton, J. D. M.-W. C., & Toutanova, L. K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT* (pp. 4171–4186).
- Koolagudi, S. G., & Rao, K. S. (2012). Emotion recognition from speech: a review. *International Journal of Speech Technology*, 15(2), 99–117.
- Lee, Y. M., Chen, Y., Li, S., & Huang, C.-r. (2010). Emotion cause events: Corpus construction and analysis. In *7th international conference on language resources and evaluation, LREC 2010* (pp. 1121–1128). European Language Resources Association (ELRA).
- Lee, J., & Hong, I. B. (2016). Predicting positive user responses to social media advertising: The roles of emotional appeal, informativeness, and creativity. *International Journal of Information Management*, 36(3), 360–373.
- Lerner, J. S., Li, Y., Valdesolo, P., & Kassam, K. S. (2015). Emotion and decision making. *Annual Review of Psychology*, 66(1), 799–823.
- Li, X., Feng, S., Wang, D., & Zhang, Y. (2019). Context-aware emotion cause analysis with multi-attention-based neural network. *Knowledge-Based Systems*, 174, 205–218.
- Li, W., Li, Y., Pandelea, V., Ge, M., Zhu, L., & Cambria, E. (2022). ECPEC: Emotion-cause pair extraction in conversations. *IEEE Transactions on Affective Computing*, 14(3), 1754–1765.
- Li, J., Lin, Z., Fu, P., & Wang, W. (2021). Past, present, and future: Conversational emotion recognition through structural modeling of psychological knowledge. In *Findings of the association for computational linguistics: EMNLP 2021* (pp. 1204–1214).

- Li, J., Meng, F., Lin, Z., Liu, R., Fu, P., Cao, Y., et al. (2022). Neutral utterances are also causes: Enhancing conversational causal emotion entailment with social commonsense knowledge. In L. D. Raedt (Ed.), *Proceedings of the thirty-first international joint conference on artificial intelligence, IJCAI-22* (pp. 4209–4215). International Joint Conferences on Artificial Intelligence Organization, <http://dx.doi.org/10.24963/ijcai.2022/584>, Main Track.
- Li, W., Shao, W., Ji, S., & Cambria, E. (2022). BiERU: Bidirectional emotional recurrent unit for conversational sentiment analysis. *Neurocomputing*, 467, 73–82.
- Li, Y., Su, H., Shen, X., Li, W., Cao, Z., & Niu, S. (2017). DailyDialog: A manually labelled multi-turn dialogue dataset. In *Proceedings of the eighth international joint conference on natural language processing (volume 1: long papers)* (pp. 986–995).
- Li, W., & Xu, H. (2014). Text-based emotion classification using emotion cause extraction. *Expert Systems with Applications*, 41(4), 1742–1749.
- Li, W., Zhu, L., Mao, R., & Cambria, E. (2023). SKIER: A symbolic knowledge integrated model for conversational emotion recognition. vol. 37, In *Proceedings of the AAAI conference on artificial intelligence* (pp. 13121–13129).
- Liu, D., Fu, J., Qu, Q., & Lv, J. (2019). Bfgan: backward and forward generative adversarial networks for lexically constrained sentence generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(12), 2350–2361.
- Liu, H., & Singh, P. (2004). ConceptNet—a practical commonsense reasoning tool-kit. *BT Technology Journal*, 22(4), 211–226.
- Liu, H., Wei, R., Tu, G., Lin, J., Liu, C., & Jiang, D. (2024). Sarcasm driven by sentiment: A sentiment-aware hierarchical fusion network for multimodal sarcasm detection. *Information Fusion*, 108, Article 102353.
- Liu, H., Yang, B., & Yu, Z. (2024). A multi-view interactive approach for multimodal sarcasm detection in social internet of things with knowledge enhancement. *Applied Sciences*, 14(5).
- Luo, M., Fei, H., Li, B., Wu, S., Liu, Q., Poria, S., et al. (2024). PanoSent: A panoptic sextuple extraction benchmark for multimodal conversational aspect-based sentiment analysis. In *Proceedings of ACM multimedia* (pp. 7667–7676).
- Majumder, N., Poria, S., Hazarika, D., Mihalcea, R., Gelbukh, A., & Cambria, E. (2019). DialogueRNN: An attentive RNN for emotion detection in conversations. In *Proceedings of the thirty-third AAAI conference on artificial intelligence and thirty-first innovative applications of artificial intelligence conference and ninth AAAI symposium on educational advances in artificial intelligence* (pp. 6818–6825).
- Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing* (pp. 1532–1543).
- Pessoa, L. (2008). On the relationship between emotion and cognition. *Nature Reviews. Neuroscience*, 9(2), 148–158.
- Picard, R. W. (2010). Affective computing: from laughter to IEEE. *IEEE Transactions on Affective Computing*, 1(1), 11–17.
- Poria, S., Majumder, N., Hazarika, D., Ghosal, D., Bhardwaj, R., Jian, S. Y. B., et al. (2021). Recognizing emotion cause in conversations. *Cognitive Computation*, 13(5), 1317–1332.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “why should i trust you?” explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1135–1144).
- Saxena, P., Huang, Y. J., & Kurohashi, S. (2022). Static and dynamic speaker modeling based on graph neural network for emotion recognition in conversation. In *Proceedings of the 2022 conference of the North American chapter of the association for computational linguistics: human language technologies: student research workshop* (pp. 247–253).
- Shen, W., Wu, S., Yang, Y., & Quan, X. (2021). Directed acyclic graph network for conversational emotion recognition. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing (volume 1: long papers)* (pp. 1551–1560).
- Sun, Y., Yen, G. G., Xue, B., Zhang, M., & Lv, J. (2021). Arctext: A unified text approach to describing convolutional neural network architectures. *IEEE Transactions on Artificial Intelligence*, 3(4), 526–540.
- Susanto, Y., Livingstone, A. G., Ng, B. C., & Cambria, E. (2020). The Hourglass model revisited. *IEEE Intelligent Systems*, 35(5), 96–102.
- Tang, C., He, Z., Li, Y., & Lv, J. (2021). Zero-shot learning via structure-aligned generative adversarial network. *IEEE Transactions on Neural Networks and Learning Systems*, 33(11), 6749–6762.
- Tu, G., Jing, R., Liang, B., Yang, M., Wong, K.-F., & Xu, R. (2023). A training-free debiasing framework with counterfactual reasoning for conversational emotion detection. In H. Bouamor, J. Pino, & K. Bali (Eds.), *Proceedings of the 2023 conference on empirical methods in natural language processing* (pp. 15639–15650). Singapore: Association for Computational Linguistics, <http://dx.doi.org/10.18653/v1/2023.emnlp-main.967>, URL: <https://aclanthology.org/2023.emnlp-main.967>.
- Tu, G., Liang, B., Jiang, D., & Xu, R. (2023). Sentiment- emotion- and context-guided knowledge selection framework for emotion recognition in conversations. *IEEE Transactions on Affective Computing*, 14(3), 1803–1816.
- Tu, G., Liang, B., Mao, R., Yang, M., & Xu, R. (2023). Context or knowledge is not always necessary: A contrastive learning framework for emotion recognition in conversations. In *Findings of the association for computational linguistics: ACL 2023* (pp. 14054–14067). Toronto, Canada: Association for Computational Linguistics, <http://dx.doi.org/10.18653/v1/2023.findings-acl.883>, URL: <https://aclanthology.org/2023.findings-acl.883>.
- Tu, G., Niu, T., Xu, R., Liang, B. B., & Cambria, E. (2024). AdaCLF: An adaptive curriculum learning framework for emotional support conversation. *IEEE Intelligent Systems*, 39(4), 5–11.
- Tu, G., Wen, J., Liu, H., Chen, S., Zheng, L., & Jiang, D. (2022). Exploration meets exploitation: Multitask learning for emotion recognition based on discrete and dimensional models. *Knowledge-Based Systems*, 235, Article 107598.
- Wang, B., Dong, G., Zhao, Y., Li, R., Cao, Q., Hu, K., et al. (2023). Hierarchically stacked graph convolution for emotion recognition in conversation. *Knowledge-Based Systems*, Article 110285.
- Wang, L., Yu, Z., Yu, H., Liu, S., Xie, Y., Guo, B., et al. (2023). AdaEvo: Edge-assisted continuous and timely DNN model evolution for mobile devices. *IEEE Transactions on Mobile Computing*.
- Xia, R., & Ding, Z. (2019). Emotion-cause pair extraction: A new task to emotion analysis in texts. In *Proceedings of the 57th annual meeting of the association for computational linguistics* (pp. 1003–1012).
- Yang, L., Shen, Y., Mao, Y., & Cai, L. (2022). Hybrid curriculum learning for emotion recognition in conversation. vol. 36, In *Proceedings of the AAAI conference on artificial intelligence* (pp. 11595–11603).
- Zhang, D., Chen, F., & Chen, X. (2023). DualGATs: Dual graph attention networks for emotion recognition in conversations. In A. Rogers, J. Boyd-Graber, & N. Okazaki (Eds.), *Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers)* (pp. 7395–7408). Toronto, Canada: Association for Computational Linguistics, <http://dx.doi.org/10.18653/v1/2023.acl-long.408>, URL: <https://aclanthology.org/2023.acl-long.408>.
- Zhang, C., Xu, Y., & Dong, Z. Y. (2020). Robustly coordinated operation of a multi-energy micro-grid in grid-connected and islanded modes under uncertainties. *IEEE Transactions on Sustainable Energy*, 11(2), 640–651.
- Zhang, D., Yang, Z., Meng, F., Chen, X., & Zhou, J. (2022). TSAM: A two-stream attention model for causal emotion entailment. In *Proceedings of the 29th international conference on computational linguistics* (pp. 6762–6772). Gyeongju, Republic of Korea: International Committee on Computational Linguistics, URL: <https://aclanthology.org/2022.coling-1.588>.
- Zhao, W., Zhao, Y., Li, Z., & Qin, B. (2023). Knowledge-bridged causal interaction network for causal emotion entailment. In *Proceedings of the thirty-seventh AAAI conference on artificial intelligence and thirty-fifth conference on innovative applications of artificial intelligence and thirteenth symposium on educational advances in artificial intelligence* (pp. 14020–14028).
- Zhong, P., Wang, D., & Miao, C. (2019). Knowledge-enriched transformer for emotion detection in textual conversations. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing* (pp. 165–176).
- Zhu, L., Mao, R., Cambria, E., & Jansen, B. J. (2024). Neurosymbolic AI for personalized sentiment analysis. In *Proceedings of the International Conference on Human-Computer Interaction* (pp. 269–290). Washington DC, USA.
- Zhu, L., Pergola, G., Gui, L., Zhou, D., & He, Y. (2021). Topic-driven and knowledge-aware transformer for dialogue emotion detection. In *Proceedings of ACL-IJCNLP* (pp. 1571–1582).

**Hao Liu** has graduated from the Department of Computer Science at Shantou University. Now he is pursuing a doctoral degree at Northwestern Polytechnical University, China. His current research interests include affective computing and deep learning.

**Runguo Wei** is currently pursuing the master's degree with the Department of Computer Science at Shantou University, China. His current research focuses on affective computing and machine learning, etc.

**Geng Tu** has graduated from the Department of Computer Science at Shantou University. His current research interests include affective computing and deep learning.

**Jiali lin** is working with the Business School of Shantou University, Shantou University, China, where she is an associate Professor.

**Dazhi Jiang** received his BA in Computer Science from the China University of Geoscience (Wuhan) in 2004. He obtained his Ph.D. from the State Key Laboratory of Software Engineering, Wuhan University, China in 2009. Since then, he has been with the Department of Computer Science, Shantou University, China, where he was a Professor. His research interests include affective computing, deep learning, data mining and applications of artificial intelligence.

**Erik Cambria** is a Professor at Nanyang Technological University, where he also holds the appointment of Provost Chair in Computer Science and Engineering. His research focuses on neurosymbolic AI for interpretable, trustworthy, and explainable affective computing in domains like social media monitoring, financial forecasting, and AI for social good. He is ranked in Clarivate's Highly Cited Researchers List of World's Top 1% Scientists, is recipient of many awards, e.g., IEEE Outstanding Early Career, was listed among the AI's 10 to Watch, and was featured in Forbes as one of the 5 People Building Our AI Future.