# Multimodal Prompt Transformer with Hybrid Contrastive Learning for Emotion Recognition in Conversation

Shihao Zou
Chongqing University of Technology
Chongqing, China
z_sh9904@163.com

Xianying Huang*
Chongqing University of Technology
Chongqing, China
wldsj_cqut@163.com

Xudong Shen
Chongqing University of Technology
Chongqing, China
sxd_cqut@163.com

## ABSTRACT

Emotion Recognition in Conversation (ERC) plays an important role in driving the development of human-machine interaction. Emotions can exist in multiple modalities, and multimodal ERC mainly faces two problems: (1) the noise problem in the cross-modal information fusion process, and (2) the prediction problem of less sample emotion labels that are semantically similar but different categories. To address these issues and fully utilize the features of each modality, we adopted the following strategies: first, deep emotion cues extraction was performed on modalities with strong representation ability, and feature filters were designed as multimodal prompt information for modalities with weak representation ability. Then, we designed a Multimodal Prompt Transformer (MPT) to perform cross-modal information fusion. MPT embeds multimodal fusion information into each attention layer of the Transformer, allowing prompt information to participate in encoding textual features and being fused with multi-level textual information to obtain better multimodal fusion features. Finally, we used the Hybrid Contrastive Learning (HCL) strategy to optimize the model's ability to handle labels with few samples. This strategy uses unsupervised contrastive learning to improve the representation ability of multimodal fusion and supervised contrastive learning to mine the information of labels with few samples. Experimental results show that our proposed model outperforms state-of-the-art models in ERC on two benchmark datasets.

## CCS CONCEPTS

• **Information systems** → **Sentiment analysis**; • **Computing methodologies** → **Natural language processing**.

## KEYWORDS

emotion recognition in conversation, multimodal prompt information, transformer, hybrid contrastive learning

*Corresponding author

## 1 INTRODUCTION

With the rapid development of social networks, there has been a lot of attention given to building dialogue systems that can understand user emotions and intentions and engage in effective dialogue interaction. Emotion Recognition in Conversation (ERC) is a task that assigns emotional labels with contextual relationships to each utterance made by speakers during a conversation. As a relevant task for dialogue systems, ERC has made important contributions to the development of engaging, interactive, and empathetic dialogue systems by analyzing user emotions in the context of a conversation. It has greatly propelled the advancement of human-machine interaction [28].
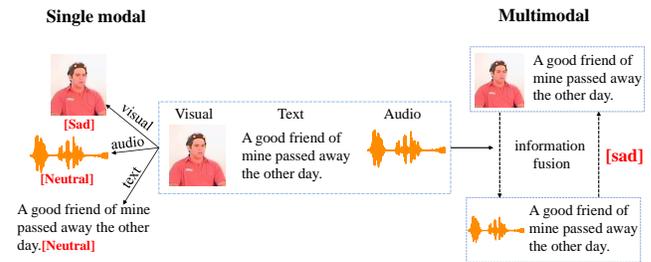


**Figure 1: Example of unimodal vs. multimodal in IEMOCAP.**

In the development of ERC, early ERC studies [8, 17, 29] focused mainly on research methods that only used text features. However, researchers found that text contains many types of utterances that are difficult for models to fully understand, such as irony, and that emotions in conversation also change with the dynamic interactive process. Since it is difficult to fully understand the conversational context, using only text is insufficient to capture emotional cues in the conversation for accurate emotion prediction. Therefore, people have attempted to expand the model's input from a single modality to multiple modalities (such as visual, audio, videos, etc.). As shown in Figure 1, using a single modality always leads to prediction errors, but using more types of modality information, the problem of semantic limitations of using only text can be mitigated. However, after introducing multiple modalities, since each modality is in a different semantic space, cross-modal information interactions will introduce a lot of noise if the semantic gaps in them are not considered. For example, Hu et al. [15] constructed a modality interaction graph by treating different modality features as nodes. Although

this method uses multiple modality information, it does not consider the semantic gap between modalities and directly performs information interaction across modalities, leading to suboptimal fusion effects and affecting the final emotion prediction.

In addition, because the ERC datasets contained many semantically similar but few sample labels (e.g., fear and disgust, happy and excited), some previous studies [13] concluded that the number of these label samples was less and their results were not statistically significant, so they were merged into other similar emotion category samples. Although this approach can improve the accuracy of emotion prediction, from a psychological perspective [1, 7, 20] , each emotion category reflects its independent emotional intensity, and simply combining emotion categories leads to oversimplification and may not accurately capture the emotional complexity and diversity of experience. Furthermore, since each emotion label represents the user's emotional state, accurate predictions should be obtained for each category label in real-world scenarios.

To address the above problems, we propose a new model for ERC, namely Multimodal Prompt Transformer with Hybrid Contrastive Learning (MPT-HCL). Firstly, in order to extract emotional cues from the text, we construct exclusive relationship graphs from both the speaker and context levels to extract emotional cues at different levels. For the audio and visual modalities, we filter the features using a designed modal feature filter, which filters out low-level features with more noise and retains high-level features with valid information. Due to previous research demonstrating the importance of text modality [39, 42], we consider text as the main modality feature and refer to the filtered features of audio and visual as textual prompt information. Then, we use MPT for information interaction between modalities. MPT embeds the interaction of multiple modalities into each attention layer of the Transformer, allowing the multimodal prompt information to participate in text feature encoding and fusion with multi-level text information, thus reducing noise generation and using the multimodal fusion result for emotion prediction of each utterance. Finally, to better optimize the model's performance, on one hand, we use unsupervised contrastive learning (UCL) to repeatedly extract mutual information [2] between the fusion features and each unimodal modality, in order to mine the relationship between modalities and optimize the fusion feature representation. On the other hand, we use supervised contrastive learning (SCL) to mine the relationship between fused features and labels in the sample. As SCL aggregates less sample labels by treating all samples with the same label in a batch as positive examples, it enhances the presence of less sample labels in the batch. Through this hybrid contrastive learning method, the feature representation of multimodal fusion can be optimized, thus effectively improving the accuracy of prediction for less sample labels. In summary, the main contributions of this paper are summarized as follows:

- We propose a novel approach of using filtered modality information as multimodal prompt information, and designing a multimodal prompt transformer for cross-modal information interaction to enhance the fusion effect of multiple modalities.
- For the first time in multimodal ERC, we introduced hybrid contrastive learning to separately explore the information between the fused modal features and each modality, as well as the information in the labels of the samples.

- We propose a new ERC model, MPT-HCL, which adopts a multimodal fusion method with hybrid contrastive learning to improve context understanding and the accuracy of multimodal ERC.
- We conducted extensive experiments on two public benchmark multimodal datasets, including IEMOCAP [3] and MELD [32]. The results showed that our proposed MPT-HCL model is more effective and superior to all SOTA baseline models.

## 2 RELATED WORKS

### 2.1 Emotion Recognition in Conversation

Emotion recognition in conversation, as an important research area in natural language processing (NLP), has received extensive attention in recent years. Existing research on ERC mainly has two types of data input, text-based and multimodal-based: (1) Text-based: DialogueGCN [8] uses graph networks for modeling dependencies between self- and inter- of speakers, which effectively solves the DialogueRNN [29] suffers from the context propagation problem; Ishiwatari et al. [17] proposes that R-GAT with relational location encoding not only captures the dependency relations between speakers, but also provides sequential information about the relational graph structure; Shen et al. [34] designs a directed acyclic graph (DAG) neural network to encode the utterance to better model the intrinsic structure in the conversation and thus explicitly model the information of each speaker in the conversation; TODKAT [41] utilizes the encoder-decoder architecture, which combines the representation of topic information with common-sense information in ERC; (2) Multimodal-based: ICON [10] and CMN [11] both model information in conversation sequences by GRU; MulT [35] uses Transformer's [37] fusion approach of the basic module-multihead attention mechanism to achieve cross-modal information fusion by using different modalities as query, key, and value in attention respectively; Li et al. [23] proposes a new structure called Emoformer to extract multimodal emotion vectors from different modalities and fuse them with sentence vectors into an emotion capsule; MM-DFN [13] uses a new multimodal dynamic fusion network to capture dynamic changes of contextual information in different semantic spaces.

### 2.2 Contrastive Learning

In the field of computer vision, SimCLR [4] optimizes contrast loss by using images obtained from the same image by randomly different data enhancement as positive samples and other images as negative samples. In natural language pre-training, ConSERT [38] introduces self-supervised contrast loss in the fine-tuning phase of BERT [5] in order to address the poor performance of sentence representation in semantic similarity tasks; Li et al. [22] uses supervised contrast learning on top of BART [21] as the backbone network to make different emotions to be mutually exclusive to better identify similar emotions. In terms of multimodal learning, TupleInfoNCE [25] is a method for learning representations of multimodal data using contrast loss that learns complementary synergies between modalities; MMIM [9] maintains task-relevant information by maximizing mutual information in single-peaked input pairs.

# 3 TASK DEFINITION

In ERC, the data consists of multiple conversations $\{c_1, c_2, \ldots, c_N\}$, and each conversation consists of a series of utterances $c_i = [u_1, u_2, \ldots, u_m]$, where $N$ is the number of conversations in a batch of data, and $m$ is the number of utterances in the i-th conversation. Each utterance $u_i$ consists of $n_i$ tokens, i.e., $\{w_{i1}, w_{i2}, \ldots, w_{in_i}\}$. Each conversation has $M$ speakers $P = \{p_1, p_2, \ldots, p_M\}, (M \geq 2)$ and each utterance is spoken by a speaker $p_{s(u_i)}$, where the function $s(\cdot)$ maps the index of the utterance to the corresponding speaker. The discrete value $y_i \in S$ is used to represent the emotion labels of $u_i$, where $S$ is the set of emotion labels. The purpose of ERC is to input a conversation and identify the correct emotion classification for each utterance in the conversation from the set of emotion labels. For each utterance $u_i$, we extract the multimodal features $u_i = \{u_i^m\}, m \in \{a, v, t\}$. Here, $u_i^a \in \mathbb{R}^{d_a}, u_i^v \in \mathbb{R}^{d_v}$ and $u_i^t \in \mathbb{R}^{d_t}$ are the audio, visual and text feature representations of the utterance, respectively. And $\{d_m\}, m \in \{a, v, t\}$ is the feature dimension of each modality.

# 4 MULTIMODAL PROMPT TRANSFORMER WITH HYBRID CONTRASTIVE LEARNING

In this section, we first introduce the learning methods for different modalities. Then, we present the designed cross-modal fusion method, the Multimodal Prompt Transformer (MPT). Finally, we show our hybrid contrastive learning (HCL) optimization strategy. The architecture of our model MPT-HCL is shown in Figure 2.

## 4.1 Contextual Capture

For each modality, we use Bi-LSTM for contextual capture, which is calculated as follows:

$$
\begin{aligned}
h_i^m &= BiLSTM^m(u_i^m, h_{i-1}^m) \\
H_m &= \{h_i^m\}_{i=1}^L \in \mathbb{R}^{L \times d_m}, m \in \{a, v, t\}
\end{aligned}
\tag{1}
$$

where $h_i^m$ is the i-th hidden layer state of Bi-LSTM. $L$ denotes the conversation sequence length in the batch, and we use $H_t = \{h_i^t\}_{i=1}^L$ as the initial representation of the node.

## 4.2 Modal Feature Filter

To reduce noise in the modal features, we have designed a modal feature filter (MFF) that selects the features with more useful information while filtering out the low-level features with excessive noise. This allows us to obtain the multimodal prompt information. The feature filter is mainly composed of a dynamic gating module. The dynamic gate predicts a normalized vector for each modality, which represents the degree to which information needs to be obtained from that modality. We take the visual modality as an example, in the dynamic gate, $z_v^{(l)}$ represents the vector that indicates the degree to which the visual modality provides information to the text modality in the l-th layer of the Transformer. We first calculate the logit of the gating signal $\theta_v$:

$$
\theta_v = f(W_l(\frac{1}{L}\sum_{i=1}^L P(h_i^v)))
\tag{2}
$$

where $f(\cdot)$ denotes the Leaky_ReLU activation function, $P(\cdot)$ represents the average pooling, which is used to generate the average

vector by weighting the average of the utterances in the current batch, and $W_l$ is the parameter matrix of the linear layer. After that, a probability vector $z_v$ is generated for the visual feature representation as follows:

$$
z_v = softmax(\theta_v)
\tag{3}
$$

Based on the probability results obtained from the dynamic gating, we obtain the final aggregated visual gating feature $V_{gate}$, and then feed the aggregation result into the lower and upper projection layers for better information fusion:

$$
\begin{aligned}
V_{gate} &= z_v H_v \\
V_d &= \sigma(W_d V_{gate} + b_d) \\
S_v &= W_u V_d + b_u
\end{aligned}
\tag{4}
$$

where $S_v \in \mathbb{R}^{L \times d}$ is the visual prompt feature representation obtained by the modal feature filter, the audio prompt information $S_a \in \mathbb{R}^{L \times d}$ is obtained in the same way as the visual, $d$ is the dimension of the modal features after alignment. $\sigma(\cdot)$ is the sigmoid activation function, and $\{W_d, W_u, b_d, b_u\}$ are the trainable parameters.
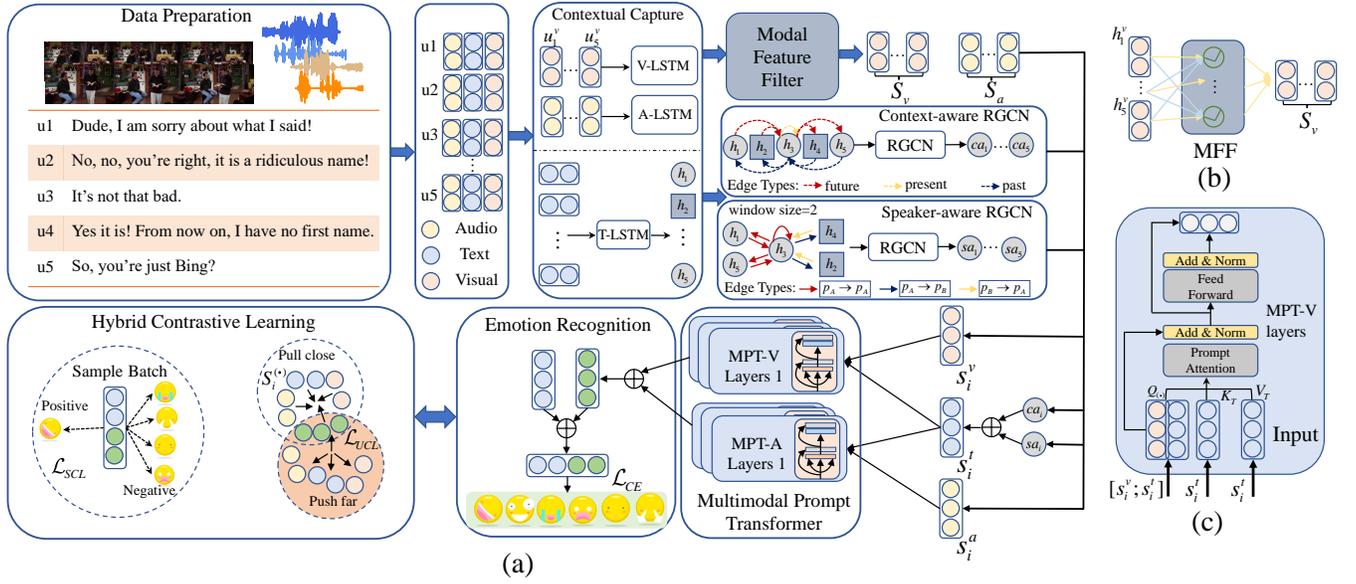
## 4.3 Speaker & Context-aware RGCN

Inspired by [24], but we don't use it for the same purpose. We design speaker-aware RGCN (Sa-RGCN) and context-aware RGCN (Ca-RGCN) as the core modules for emotion cue extraction in text modality. In these modules, edges are used to measure the importance of connections between nodes. The type of edge defines the method of propagation of different information between nodes. Sa-RGCN and Ca-RGCN have the same edges, but each edge represents a different dependency, and we will describe the composition of these two modules in detail next.

**Edge:** For each node, its interaction with the context nodes should be considered. If each node $v_i$ interacts with all the context nodes $v_j$, the constructed graph contains a large number of edges, which usually leads to the problem of computational difficulties due to the huge number of parameters. To solve this problem, we fix the context window size to $w$, so that each node $v_i$ interacts with only $\{v_j\}_{j=max(i-w,1)}^{min(i+w,L)}$ context nodes. In our model implementation, we perform the selection of $w$ within $w \in \{1, 2, 3, 4\}$ and the edges $e_{ij}$ denote nodes $v_i$ to $v_j$.

**Speaker-aware RGCN:** Sa-RGCN uses different speakers and their spoken utterances to capture the dependencies of the speakers in a conversation. Specifically, we assign a speaker identifier $\alpha_{ij} \in \alpha$ to each edge $e_{ij}$. Here $\alpha$ denotes the set of speaker types in the conversation, and $|\alpha|$ denotes the number of $\alpha$. For each edge $e_{ij}$, $\alpha_{ij}$ serves as the set of $p_{s(u_i)} \rightarrow p_{s(u_j)}$, where $p_{s(u_i)}$ and $p_{s(u_j)}$ denote the speaker identifiers of $u_i$ and $u_j$, respectively.

**Context-aware RGCN:** Ca-RGCN uses contextual information to capture the contextual dependencies in a conversation. Specifically, we assign a context type identifier $\beta_{ij} \in |\beta|$ to each edge $e_{ij}$, where $\beta$ denotes the set of context types in the conversation. Based on the relative positions of $u_i$ and $u_j$ in the conversation, we will perform $\beta_{ij}$ value determination from $\{past, present, future\}$, so $|\beta|=3$.

**Graph learning:** We use RGCN to aggregate the neighbor information in the graph. For Sa-RGCN and Ca-RGCN, we pass different

Figure 2: (a) shows the specific structure of MPT-HCL. (b) illustrates the architecture of the Modality Feature Filter (MFF) module, showcasing its input and output representations. (c) provides a detailed view of MPT-V, designed specifically for the visual modality.

information through the edges, and the parameters depend on the type of the edges. The calculation is shown as follows:

$$sa_i = ReLU(\sum_{r \in \alpha} \sum_{j \in N_i^r} \frac{1}{|N_i^r|} W_r^s h_j^t)$$

$$ca_i = ReLU(\sum_{r \in \beta} \sum_{j \in N_i^r} \frac{1}{|N_i^r|} W_r^c h_j^t) \tag{5}$$

where $sa_i$ and $ca_i$ denote the outputs of nodes in Sa-RGCN and Ca-RGCN, respectively. $N_i^r$ denotes the set of all neighbor nodes of $v_i$ under relation $r$, and $|N_i^r|$ is the number of $N_i^r$. $W_r^s$ and $W_r^c$ denote the parameter matrices of Sa-RGCN and Ca-RGCN at relation $r$, respectively.

After obtaining the speaker and context dependencies in the text, we fuse them to obtain the enhanced textual feature representation $s_i^t$, which is computed as follows:

$$s_i^t = sa_i + ca_i \tag{6}$$

## 4.4 Multimodal Interaction

We consider the audio and visual features obtained by the modal feature filters module as textual prompt information, and update all text states with the prompt information when executing cross-modal attention sequentially. In this way, the final fusion feature will have the ability to encode both the context and cross-modal semantic information, which effectively alleviates the problem of noise generated by unrelated elements in the process of multimodal fusion.

**Prompt Attention.** First, we project the textual feature sequence $S_t = \{s_i^t\}_{i=1}^L \in \mathbb{R}^{L \times d}$ enhanced with emotional cues into the query/key/value vector, i.e., $Q = S_t W_Q, K = S_t W_K, V = S_t W_V$,

where $W_{(.)}$ is the trainable parameter matrix. Then we prepend the visual prompt feature sequences $S_v = \{s_i^v\}_{i=1}^L \in \mathbb{R}^{L \times d}$ and audio prompt feature sequences $S_a = \{s_i^a\}_{i=1}^L \in \mathbb{R}^{L \times d}$ obtained from the modal feature filters to each attention layer of the Transformer, respectively, and we call this proposed prepending method Prompt Attention (PA), by which the prompt features are used for effective cross-modal interaction, as implemented in the following equation shown:

$$Y_{tv} = PA([S_v; S_t], S_t, S_t)$$

$$= softmax(\frac{[S_v; Q]K^T}{\sqrt{d}})V \tag{7}$$

where $[;]$ denote feature concatenate, $Y_{tv} \in \mathbb{R}^{L \times d}$ is the result of the weight of a layer of PA. In this approach, multiple attentions are combined to obtain the output results of the multihead attention layer as follows:

$$Multihead(Y_{tv}) = Concat(Y_{tv}^1, \ldots, Y_{tv}^n)W \tag{8}$$

where $Y_{tv}^1, \ldots, Y_{tv}^n$ is the output of each attention layer, $n$ is the number of attention layers, and $W$ is the trainable parameter matrix. PA is used to perform cross-modal interaction processes between text modalities and multimodal prompt information.

**Multimodal Prompt Transformer.** Based on the above PA, we design the Multimodal Prompt Transformer (MPT) with the structure shown in Figure 2(c). Since we embed multimodal interactions into each attention layer of the Transformer, visual and audio features can participate in the encoding of text features and fuse with multi-layered text information. Low-level syntactic features encoded by the shallow Transformer layer and high-level semantic features encoded by the deep Transformer layer [30, 37] interact with the visual and audio prompt features, enabling the

fusion of information across modalities. In addition depending on the multimodal prompt features used, as more effective feature representation can fully exploit the fusion effect between modalities.

We use a residual connectivity layer with regularization to normalize the output of the multihead attention layer of Eq.(8) and use a position feedforward sublayer to obtain the output of the attention:

$$
\begin{aligned}
N &= Norm(Y_{tv} + Multihead(Y_{tv})) \\
F &= max(0, NW_1 + b_1)W_2 + b_2 \\
G &= Norm(F + Multihead(F))
\end{aligned}
\tag{9}
$$

where $W_1, W_2$ is the weight parameter and $b_1, b_2$ is the bias parameter. In this process each modal prompt feature is continuously updated with the feature of the text by the above method to obtain better fusion results, and finally we use self-attention to collect the sequence information of the modal fusion features to obtain the multimodal fusion result $X \in \mathbb{R}^{L \times d}$ for the current conversation:

$$
X = Attention(G) \tag{10}
$$

Connecting the above Eq.(8) to (10), we can input different modal prompt features and text features to MPT to obtain the modal fusion feature representation:

$$
\begin{aligned}
X_{tv} &= MPT([S_v; S_t], S_t, S_t) \\
X_{ta} &= MPT([S_a; S_t], S_t, S_t)
\end{aligned}
\tag{11}
$$

where $S_t, S_v, S_a$ denote the text, visual and audio features obtained by different learning methods as the input to the MPT, and $X_{tv}, X_{ta}$ denote the multimodal fusion results obtained after prompting the text with visual prompt features and audio prompt features, respectively. We concatenate the outputs from different MPTs to obtain the fused features $X_{mpt} \in \mathbb{R}^{L \times 2d}$ :

$$
X_{mpt} = X_{tv} \oplus X_{ta} \tag{12}
$$

Finally, in order to fully utilize the text feature representation, we combine $X_{mpt}$ with the text modality to obtain the final representation of the current utterance used for emotion prediction $X_{fusion}$:

$$
X_{fusion} = S_t \oplus X_{mpt} \tag{13}
$$

## 4.5 Hybrid Contrastive Learning

**Unsupervised contrastive learning.** Although a better representation of the fusion features is obtained by the MPT, the relationship between each unimodal modal and the fusion features is not entirely explored, so we use unsupervised contrastive learning(UCL) to exploit the connection between them and thus optimize the obtained fusion features. We repeat the mutual information maximization between the fusion results and the input modalities, and the optimization goal is to fuse the network from each unimodal modal to the fusion features. Since we now obtain the multimodal fusion result $X_{mpt}$ by the constructed MPT network , but the mining of the connection from the fusion feature $X_{mpt}$ to each unimodal input $S_x, x \in \{a, v, t\}$ is missing. So we follow the operation of [36] and use the score function $Score(\cdot)$ with normalized prediction and true vector to measure the connection between them as follows:

$$
Score(s_x, X_{mpt}) = exp(\bar{s}_x(\overline{G}_\varphi(X_{mpt}))^T)
$$
$$
\overline{G}_\varphi(X_{mpt}) = \frac{G_\varphi(X_{mpt})}{||G_\varphi(X_{mpt})||_2}, \bar{s}_x = \frac{s_x}{||s_x||_2}
\tag{14}
$$

where $G_\varphi$ is a neural network with parameter $\varphi$ that generates a prediction for $s_x$ from $X_{mpt}$, and $|| \cdot ||_2$ is the Euclidean norm. We treat all other representations of the modality in the same batch as negative samples, and thus calculate the loss between the individual modality and the fused features:

$$
\mathcal{L}(X_{mpt}, S_x) = -\mathbb{E}_s\left[log\frac{Score(X_{mpt}, s_x^i)}{\sum_{s_x^j \in S_x} Score(X_{mpt}, s_x^j)}\right] \tag{15}
$$

Finally, the loss function of UCL consists of the losses between the fused features $X_{mpt}$ (noted here as $x$) and the text, visual and audio, respectively:

$$
\mathcal{L}_{UCL} = \mathcal{L}^{x,v} + \mathcal{L}^{x,a} + \mathcal{L}^{x,t} \tag{16}
$$

**Supervised contrastive learning.** Supervised contrastive learning(SCL) assumes that attention will be paid to certain key labels, and then treats all samples in the batch with the same label as positive samples and those with different labels as negative samples by making full use of the label information. Specifically for ERC, since the number of samples in each category in the dataset is very imbalanced, such that the information of these samples is masked in the process of calculating the loss. In addition, if only one sample exists in a category batch, it cannot be directly applied to the loss calculation. Therefore, in order to ensure that a sufficient number of feature representations are available for the loss calculation each time, we combine the obtained multimodal fusion features and text features at the position of sequence length. For a batch with $L$ training samples, we can obtain $2L$ samples in this way, after which the SCL loss $\mathcal{L}_{SCL}$ is calculated as follows:

$$
C = [S_t; X_{mpt}] \tag{17}
$$

$$
\mathcal{L}_{SCL} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} SIM(p, i)
$$
$$
SIM(p, i) = log\frac{exp((C_i \cdot C_p)/\tau)}{\sum_{a \in A(i)} exp(C_i \cdot C_p/\tau)}
\tag{18}
$$

where $C \in \mathbb{R}^{2L \times d}, i \in I = \{1, 2, \ldots, 2L\}$ denotes the index of samples in a batch, $\tau \in R^+$ denotes the temperature coefficient used to control the distance between samples, $P(i) = I_{j=i} - \{i\}$ denotes the samples with the same emotion category as $i$ but excluding $i$ itself, $|P(i)|$ denotes the number of samples, and $A(i) = I - \{i\}$ denotes the samples in a batch other than itself.

## 4.6 Model Training

The loss of the model training process consists of three components: the logarithmic loss of standard cross-entropy $\mathcal{L}_{CE}$ , the supervised contrastive loss $\mathcal{L}_{SCL}$ and the unsupervised contrastive loss $\mathcal{L}_{UCL}$.

$$
\begin{aligned}
P_i &= softmax(W_s X_{fusion} + b_s) \\
\hat{y}_i &= argmax(P_i)
\end{aligned}
\tag{19}
$$

$$
\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{J} y_{i,j} \cdot log\hat{y}_{i,j} \tag{20}
$$

$$
\mathcal{L} = \mathcal{L}_{CE} + \lambda_1 \mathcal{L}_{SCL} + \lambda_2 \mathcal{L}_{UCL} \tag{21}
$$

where $N$ is the number of conversations, $J$ denotes the number of emotion categories, $y_{i,j}$ is the true emotion label of utterance $i$, $\hat{y}_{i,j}$ denotes the probability distribution that the prediction of utterance

$i$ is category $j$, and $\lambda_1, \lambda_2$ denote the weights of supervised and unsupervised contrastive loss, respectively. In the training process we use the Adam [19] optimizer with stochastic gradient descent to train our network model.

## 5 EXPERIMENT

### 5.1 Datasets and Evaluations

We evaluated the effectiveness of our model on two benchmark datasets, IEMOCAP [3] and MELD [32]. Both datasets are multimodal ERC datasets containing text, audio, and visual. We divided the datasets according to [13]. The data distribution of two datasets are shown in Table 1 and emotion distribution information of two datasets are shown in Table 2.

**Table 1: Data distribution of IEMOCAP and MELD.**

| Dataset | Conversation | | Utterance | | Classes |
|---------|------------|------|-----------|------|---------|
| | Train+Val | Test | Train+Val | Test | |
| IEMOCAP | 120 | 31 | 5810 | 1623 | 6 |
| MELD | 1153 | 280 | 11098 | 2610 | 7 |

**IEMOCAP:** The multimodal ERC dataset. Each conversation in IEMOCAP is from two actors' performances based on the script. There are 7433 utterances and 151 conversations in IEMOCAP. Each utterance in the conversation is labeled with six categories of emotions: *happy, sad, neutral, angry, excited*, and *frustrated*.

**MELD:** The data were obtained from the TV show *Friends* and included a total of 13708 utterances and 1433 conversations. Unlike the IEMOCAP dyadic dataset, MELD has three or more speakers in a conversation, and each utterance in the conversation is labeled with seven categories of emotions: *neutral, surprise, fear, sadness, joy, disgust*, and *anger*.

**Evaluation Metrics:** We use the F1-score to evaluate the performance for each emotion class and use the weighted average of accuracy and F1-score(W-F1) to evaluate the overall performance on the two datasets.

### 5.2 Baselines

- **BC-LSTM** [31] encodes contextual semantic information through a Bi-LSTM network, thus making emotion predictions.
- **ICON** [10] uses two GRUs to model the speaker's information, additional global GRUs are used to track changes in emotional states throughout the conversation, and a multi-layer memory network is used to model global emotional states.
- **DialogueRNN** [29] models the speaker and sequential information in a conversation through three different GRUs (global GRU, speaker GRU, and emotion GRU).
- **DialogueGCN** [8] applies GCN to ERC, and the generated features can integrate rich information. RGCN and GCN are both nonspectral domain GCN models for encoding graphs.
- **DialogueXL** [33] uses the XLNet model for ERC to obtain global contextual information.

- **DialogueCRN** [14] introduces a cognitive phase that extracts and integrates emotional cues from the context retrieved during the perception phase.
- **BiDDIN** [40] specializes inter-modal and intra-modal interactive modules for corresponding modalities, as well as models contextual influence with an extra positional attention. It is set under the multimodal scenario and employs separate modality-shared and modality-specific modules.
- **MMGCN** [15] uses GCN networks to obtain contextual information, which can not only make use of multimodal dependencies effectively, but also leverage speaker information.
- **MVN** [27] explores the emotion representation of the query utterance from the word- and utterance-level views, which can effectively capture the word-level dependencies among utterances and utterance-level dependencies in the context.
- **CoG-BART** [22] uses supervised contrastive learning to better distinguish between similar emotional labels and augments the model's ability to handle context with an auxiliary generation task.
- **MM-DFN** [13] fuses multimodal contextual information by designing a new graph-based dynamic fusion module to fully understand multimodal conversational contexts to recognize emotions in utterances.
- **COGMEN** [18] is a multimodal context-based graph neural network which using local information ( speaker information) and global information (contextual information) in the conversation.
- **EmoCaps** [23] uses the new structure Emoformer to extract emotion vectors from different modalities and fuse them with sentence vectors into a emotion capsule for emotion prediction of utterance.

### 5.3 Implementation Details

Our proposed model is implemented on the Pytorch framework. The hyperparameters are set as follows: the number of layers of the Transformer in MPT is 5 ($l = 5$), where the number of prompt attention heads is 5 ($m = 5$). The coefficient $\lambda_1$ for supervised contrast loss in hybrid contrastive learning is 0.1, and the coefficient $\lambda_2$ for unsupervised contrast loss is 0.05. The dropout in both IEMOCAP and MELD is 0.2. The learning rate in IEMOCAP is 0.0001 and in MELD is 0.0003. Each training and testing procedure was run on a single RTX 3090 GPU, and the reports of our implemented models are based on the average scores of five random runs on the test set.

### 5.4 Multimodal Feature Extraction

In this paper, we use pre-extracted unimodal features following identical extraction procedures as previous methods [11, 14, 34].

**Text Feature.** To extract better utterance representation with strong representational ability, we use the large general pretrained language model RoBERTa-Large [26] for text vector encoding extraction. However, unlike other downstream tasks, we use the transformer structure to encode the utterances without classifying or decoding them. More specifically, for each utterance in the text modal, we precede its token with a special token $[CLS]$ to make it of the form of $\{[CLS], w_{i1}, w_{i2}, \ldots, w_{in_i}\}$. Then we use the pooled

**Table 2: The number of utternces with each emotion label in the IEMOCAP and MELD test dataset.**

| Dataset | Emotions | | | | | | | | |
|---------|---------|----------|------|------------|-----------|---------|-------|---------|------------|
| | Neutral | Surprise | Fear | Sad/Sadness | Happy/Joy | Disgust | Angry | Excited | Frustrated |
| IEMOCAP | 384 | - | - | 245 | 144 | - | 170 | 299 | 381 |
| MELD | 1256 | 281 | 50 | 208 | 402 | 68 | 345 | - | - |

embedding result of the last layer of [*CLS*] as the feature representation of $u_i^T$, and finally, we obtain a sentence vector with 1024 dimensions for each utterance.

**Audio and Visual Feature.** In terms of audio features, OpenSmile [6] is used with the IS13 comparison profile, which extracted a total of 6373 features for each utterance video, we reduced the dimensionality to 1582 for the IEMOCAP and 300 for the MELD dataset by using a fully connected layer. The visual facial features were extracted by pretraining on the Facial Expression Recognition Plus (FER+) corpus using DenseNet [16]. This captures changes in the expression of the speakers, which is very important information for ERC. Finally, a 342-dimensional visual feature representation was obtained.

## 6 RESULTS AND ANALYSIS

### 6.1 Main Result

Tables 3 and Tables 4 show the experimental results of our proposed MPT-HCL model and the baselines on the IEMOCAP and MELD datasets. Models marked with "*" only use the text modality, and "-" indicates that the results were not reported. The best results are highlighted in bold. All other results are reported in their respective papers. On the one hand, compared to existing methods, our model achieved better utterance representations through pretrained language models in sentence encoding. On the other hand, as shown in Tables 3 and Tables 4, we found that: (1) Our proposed MPT-HCL outperforms all baseline models in both evaluation metrics, demonstrating the effectiveness of our model in multimodal ERC. (2) MPT-HCL achieves the best results in almost all emotion categories when compared to EmoCaps. Taking the MELD dataset as an example, our approach outperforms EmoCaps by a large margin in predicting labels with few samples, such as Fear and Disgust. This demonstrates the effectiveness of our multimodal fusion method and the specifically designed hybrid contrastive learning strategy for handling labels with few samples. (3) The overall effect of MPT-HCL is better than MM-DFN, which is a recent model using speaker information.This highlights the effectiveness of utilizing dialogue information in our approach. We consider the reasons for emotional changes in human communication in real life, which are often influenced either by the words spoken by others or by one's own mood. Therefore, we extract contextual cues from both the speaker's own utterances(Speaker-aware) and the interaction between speakers(Context-aware), which leads to superior results compared to MM-DFN.

### 6.2 Ablation Study

*6.2.1 Various Modalities.* Table 5 shows the performance of our model on the MELD and IEMOCAP datasets under different modal combinations. We can observe that: (1) The performance of multimodal inputs is better than that of unimodal inputs. Furthermore, among the three modalities of text, audio, and visual, just as we intended to use text as the main modality and other auxiliary modality features as prompt, the text modality performs better than the other two modalities. (2) After adding the text modality, the combination of visual and audio shows a significant improvement in performance. This suggests that the text modality plays an important role in ERC, while audio and visual serve as auxiliary features to improve the accuracy of the model's recognition. This is consistent with our goal of filtering the audio and visual modalities to extract more effective features. (3) When comparing the normal modal combination (i.e., T+A+V) with our feature-filtered modal combination, we found that the latter is more effective for modality interaction. This is because the visual and audio modals were originally used as auxiliary modalities. By further filtering their features and retaining their useful high-level features as prompt information for the text modality, we obtained a more effective multimodal fusion feature representation.

*6.2.2 Module Analysis.* To study the contribution of each module in MPT-HCL, we performed an ablation study on both datasets. We consider the following settings:

- **Full A:** We do not use MFF to filter audio feature.
- **Full V:** We do not use MFF to filter visual feature.
- **w/o MPT:** We remove the multimodal prompt transformer module.
- **w/o UCL:** We remove the unsupervised contrastive learning module.
- **w/o SCL:** We remove the supervised contrastive learning module.
- **w/o Sa/Ca-RGCN:** We remove the module used to extract emotional cues from the text modal.

Table 6 shows the results of our ablation experiments, from which we can conclude that: (1) With Full A/V, we can see that unfiltered modality features introduce a lot of noise when exchanging information across modalities, and the visual modality has more noise than the audio modality due to the complex scenes in the data. (2) We can observe that not using MPT yields poor performance. This can be attributed to the lack of information interaction between modalities and the inability of simple concatenation to leverage complementary effects. It introduces significant noise and leads to a confused feature representation, resulting in subpar performance. (3) Our proposed HCL is very effective as removing any contrastive loss leads to a decrease in model performance. This module performs better on the MELD dataset, as this dataset has more minority class labels. (4) After removing the Sa/Ca-RGCN emotion cue extraction module, a significant performance drop was

**Table 3: Results of different models on the IEMOCAP dataset.**

| Model | IEMOCAP | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Happy | Sad | Neutral | Angry | Excited | Frustrated | ACC | W-F1 |
| BC-LSTM | 43.30 | 69.28 | 55.84 | 61.80 | 59.33 | 60.20 | - | 59.19 |
| CMN | 30.38 | 62.41 | 52.39 | 59.83 | 60.25 | 60.69 | - | 56.13 |
| DialogueRNN* | 33.18 | 78.80 | 59.21 | 65.28 | 71.86 | 58.91 | - | 62.75 |
| DialogueGCN* | 43.57 | 80.48 | 57.69 | 53.95 | 72.81 | 57.33 | 63.22 | 62.89 |
| DialogueXL* | - | - | - | - | - | - | - | 65.94 |
| DialogueCRN* | - | - | - | - | - | - | 66.05 | 59.19 |
| BiDDIN | - | - | - | - | - | - | 65.60 | 65.30 |
| MVN* | 55.75 | 73.30 | 61.88 | 65.96 | 69.50 | 64.21 | 65.32 | 65.44 |
| CoG-BART* | - | - | - | - | - | - | - | 66.18 |
| MMGCN | 42.34 | 78.67 | 61.73 | 69.00 | 74.33 | 62.32 | - | 66.22 |
| MM-DFN | 42.22 | 78.89 | 66.42 | 69.77 | 75.56 | 66.33 | 68.21 | 68.18 |
| COGMEN | - | - | - | - | - | - | 68.20 | 67.60 |
| EmoCaps | **71.91** | 85.06 | 64.48 | 68.99 | **78.41** | 66.76 | - | 71.77 |
| Ours(MPT-HCL) | 58.13 | **85.97** | 66.75 | **69.96** | 74.06 | **69.06** | **72.83** | **72.51** |

**Table 4: Results of different models on the MELD dataset.**

| Model | MELD | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Neutral | Surprise | Fear | Sadness | Joy | Disgust | Angry | ACC | W-F1 |
| BC-LSTM | 73.80 | 47.70 | 5.40 | 25.10 | 51.30 | 5.20 | 38.40 | 59.62 | 57.29 |
| DialogueRNN* | 73.50 | 49.40 | 1.20 | 23.80 | 50.70 | 1.70 | 41.50 | 60.31 | 57.66 |
| DialogueGCN* | - | - | - | - | - | - | - | 58.62 | 56.36 |
| DialogueXL* | - | - | - | - | - | - | - | - | 62.41 |
| DialogueCRN* | - | - | - | - | - | - | - | 61.11 | 58.67 |
| MVN* | 76.65 | 53.18 | 11.70 | 21.82 | 53.62 | 21.86 | 42.55 | 61.29 | 59.03 |
| CoG-BART* | - | - | - | - | - | - | - | - | 64.81 |
| MMGCN | - | - | - | - | - | - | - | 60.42 | 58.31 |
| MM-DFN | 77.76 | 50.69 | - | 22.93 | 54.78 | - | 47.82 | 62.49 | 59.46 |
| EmoCaps | 77.12 | **63.19** | 3.03 | 42.52 | 57.50 | 7.69 | 57.54 | - | 64.00 |
| Ours(MPT-HCL) | **77.82** | 58.26 | **21.52** | **45.15** | **60.18** | 30.36 | **59.25** | **65.86** | **65.02** |

**Table 5: Performance of MPT-HCL under different multi-modal settings.**

| | IEMOCAP | | MELD | |
|---|---|---|---|---|
| | ACC | W-F1 | ACC | W-F1 |
| Only T | 68.92 | 68.35 | 61.89 | 61.63 |
| Only A | 46.58 | 47.26 | 44.12 | 40.15 |
| Only V | 39.28 | 39.75 | 36.25 | 35.69 |
| T+A | 70.09 | 69.81 | 62.84 | 61.12 |
| T+V | 69.79 | 68.31 | 61.11 | 60.03 |
| T+A+V | 71.02 | 70.22 | 62.55 | 62.13 |
| Ours | **72.83** | **72.51** | **65.86** | **65.02** |

**Table 6: Ablation studies on various modules.**

| | IEMOCAP | | MELD | |
|---|---|---|---|---|
| | ACC | W-F1 | ACC | W-F1 |
| Full A | 71.43 | 71.08 | 64.44 | 64.17 |
| Full V | 71.13 | 69.89 | 64.11 | 64.02 |
| w/o MPT | 45.57 | 45.98 | 43.91 | 42.68 |
| w/o UCL | 71.27 | 71.11 | 63.91 | 63.62 |
| w/o UCL | 71.34 | 71.23 | 64.41 | 63.25 |
| w/o Sa/Ca-RGCN | 64.10 | 63.54 | 61.30 | 62.22 |
| Ours | **72.83** | **72.51** | **65.86** | **65.02** |

observed on both datasets. This is because both datasets involve conversation between two or more speakers, and capturing emotional cues from speakers and contextual information in conversation is crucial. However, although we constructed a speaker-aware relational graph and utilized the relationships between speakers, we did not fully explore the independent information of each speaker, especially in the case of MELD, which consists entirely of multi-party dialogues. Therefore, the performance of this module on the MELD dataset was not as significant as that on the IEMOCAP dataset.

## 6.3 The Potency of Hybrid Contrastive Learning

In order to conduct a qualitative analysis of the hybrid contrastive learning(HCL), we used t-SNE [12] to visualize the initial distribution of some data and the hidden layer status of the model after using HCL. As shown in Figure 3, when HCL loss is not used, the samples between different labels are randomly scattered, and some samples with similar emotions also overlap, which increases the difficulty of the model learning decision boundaries and leads to large errors in predicting similar labels. With the use of HCL, it can be clearly seen that the coupling degree between different categories gradually increases, resulting in a significant category aggregation effect. This shows that the hybrid contrastive learning strategy we designed plays an important role in sample classification. It is worth noting that although our hybrid contrastive learning has already achieved good category division of samples, we can still see obvious errors in some samples. We analyzed the reason to be that the number of this sample compared to the other labeled samples in the entire batch was the highest, which led to the problem of predicting other labels as this category, and this is a problem that we need to further address.
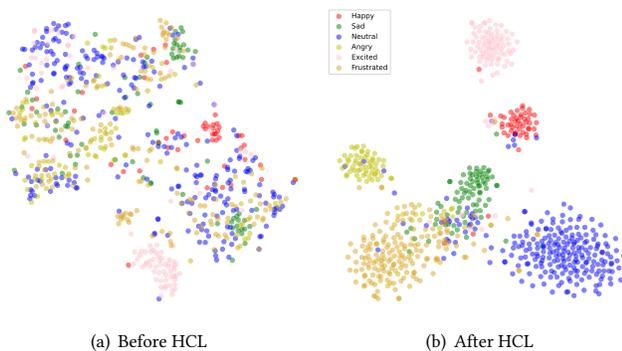


(a) Before HCL  (b) After HCL

**Figure 3: The t-SNE visualization results of the model output when HCL is use or not.**

## 6.4 Case Study

We compare the predictions of different methods for the simultaneous presence of emotion transfer as well as category less labels in utterances. Figure 4 provides a representative example of the MELD test set. In this example, $P_A$ and $P_B$ try around whether the chick can swim or not, resulting in a final finding that the chick is about to get into danger resulting in an emotion shift from neutral to fear. We observe that our model is more accurate in making emotion predictions compared to the unimodal DialogueRNN and CoG_BART models for the following reasons: for turn 4, the utterance ends up being incorrectly predicted due to its forward position without much contextual information and without other modal information as an aid. In MM-DFN, for turn 7, the utterance had too little content and the corresponding other modal features could not provide effective help, which led to the model not getting a good fusion representation in the multimodal fusion process of the sentence, resulting in the wrong prediction of the final emotion.

Our model first enhances the extracted conversational context with two levels of emotional cue extraction to avoid the problem of capturing insufficient feature content, and also adds less sample labels to the model when predicting the emotion after obtaining valid features by using a hybrid contrastive learning method. Therefore, our model achieves good results in the whole prediction process in both emotion transfer and few-sample prediction.

## 7 CONCLUSION

We present a new multimodal prompt transformer with hybrid contrastive learning (MPT-HCL) model for ERC. The RGCN is used to extract emotional clues in the conversation, and different levels of emotional clues are extracted to enhance the text modality. Modal feature filters are designed for the visual and audio modalities to filter features and obtain prompt information for the text modality, through MPT to result in better multimodal fusion feature representation. On this basis, hybrid contrastive learning is used to optimize the fusion feature representation and to explore the information of few labels in the samples, thereby improving the prediction accuracy of few sample labels. In the future, we will explore the joint methods for identifying emotional causes and emotional recognition to alleviate the problem of error propagation in information.

## REFERENCES

[1] F Gregory Ashby, Alice M Isen, et al. 1999. A neuropsychological theory of positive affect and its influence on cognition. *Psychological review* 106, 3 (1999), 529.

[2] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, R. Devon Hjelm, and Aaron C. Courville. 2018. Mutual Information Neural Estimation. In *Proceedings of the 35th International Conference on Machine Learning, ICML (Proceedings of Machine Learning Research, Vol. 80)*. 530–539. http://proceedings.mlr.press/v80/belghazi18a.html

[3] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N. Chang, Sungbok Lee, and Shrikanth S. Narayanan. 2008. IEMOCAP: interactive emotional dyadic motion capture database. *Lang. Resour. Evaluation* 42, 4 (2008), 335–359. https://doi.org/10.1007/s10579-008-9076-6

[4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. A Simple Framework for Contrastive Learning of Visual Representations. In *Proceedings of the 37th International Conference on Machine Learning, ICML (Proceedings of Machine Learning Research, Vol. 119)*. 1597–1607. http://proceedings.mlr.press/v119/chen20j.html

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, Volume 1 (Long and Short Papers)*. 4171–4186. https://doi.org/10.18653/v1/n19-1423

[6] Florian Eyben, Martin Wöllmer, and Björn W. Schuller. 2010. Opensmile: the munich versatile and fast open-source audio feature extractor. In *Proceedings of the 18th International Conference on Multimedia 2010*. ACM, 1459–1462. https://doi.org/10.1145/1873951.1874246

[7] Maria Gendron, Kristen A Lindquist, Lawrence Barsalou, and Lisa Feldman Barrett. 2012. Emotion words shape emotion percepts. *Emotion* 12, 2 (2012), 314.

[8] Deepanway Ghosal, Navonil Majumder, Soujanya Poria, Niyati Chhaya, and Alexander F. Gelbukh. 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*
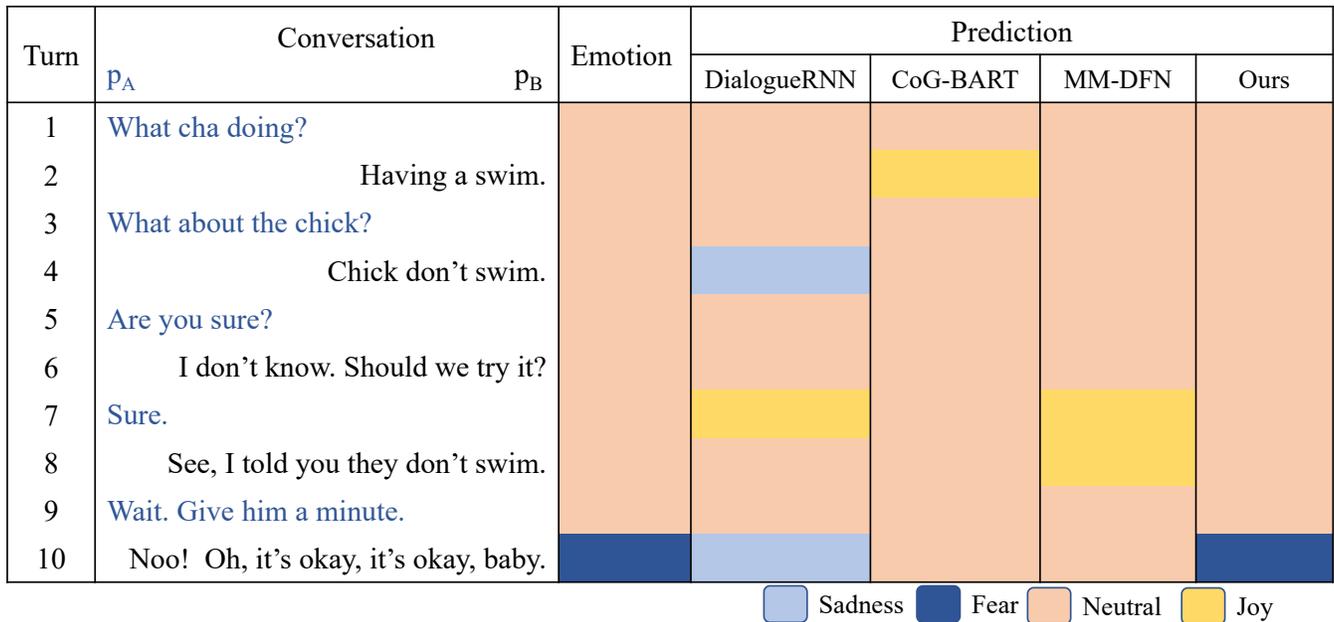
| Turn | Conversation | | Emotion | Prediction | | | |
|---|---|---|---|---|---|---|---|
| | $p_A$ | $p_B$ | | DialogueRNN | CoG-BART | MM-DFN | Ours |
| 1 | What cha doing? | | | | | | |
| 2 | | Having a swim. | | | | | |
| 3 | What about the chick? | | | | | | |
| 4 | | Chick don't swim. | | | | | |
| 5 | Are you sure? | | | | | | |
| 6 | | I don't know. Should we try it? | | | | | |
| 7 | Sure. | | | | | | |
| 8 | | See, I told you they don't swim. | | | | | |
| 9 | Wait. Give him a minute. | | | | | | |
| 10 | | Noo!  Oh, it's okay, it's okay, baby. | | | | | |

☐ Sadness   ■ Fear   ☐ Neutral   ☐ Joy

**Figure 4: Case study in MELD.**

*International Joint Conference on Natural Language Processing, EMNLP-IJCNLP.* 154–164. https://doi.org/10.18653/v1/D19-1015

[9] Wei Han, Hui Chen, and Soujanya Poria. 2021. Improving Multimodal Fusion with Hierarchical Mutual Information Maximization for Multimodal Sentiment Analysis. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP.* 9180–9192. https://doi.org/10.18653/v1/2021.emnlp-main.723

[10] Devamanyu Hazarika, Soujanya Poria, Rada Mihalcea, Erik Cambria, and Roger Zimmermann. 2018. ICON: Interactive Conversational Memory Network for Multimodal Emotion Detection. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP.* 2594–2604. https://doi.org/10.18653/v1/d18-1280

[11] Devamanyu Hazarika, Soujanya Poria, Amir Zadeh, Erik Cambria, Louis-Philippe Morency, and Roger Zimmermann. 2018. Conversational Memory Network for Emotion Recognition in Dyadic Dialogue Videos. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, Volume 1 (Long Papers).* 2122–2132. https://doi.org/10.18653/v1/n18-1193

[12] Geoffrey E Hinton and Sam Roweis. 2002. Stochastic neighbor embedding. *Advances in neural information processing systems* 15 (2002).

[13] Dou Hu, Xiaolong Hou, Lingwei Wei, Lian-Xin Jiang, and Yang Mo. 2022. MM-DFN: Multimodal Dynamic Fusion Network for Emotion Recognition in Conversations. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP.* 7037–7041. https://doi.org/10.1109/ICASSP43922.2022.9747393

[14] Dou Hu, Lingwei Wei, and Xiaoyong Huai. 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, ACL/IJCNLP, (Volume 1: Long Papers).* 7042–7052. https://doi.org/10.18653/v1/2021.acl-long.547

[15] Jingwen Hu, Yuchen Liu, Jinming Zhao, and Qin Jin. 2021. MMGCN: Multimodal Fusion via Deep Graph Convolution Network for Emotion Recognition in Conversation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP (Volume 1: Long Papers).* 5666–5675. https://doi.org/10.18653/v1/2021.acl-long.440

[16] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 4700–4708.

[17] Taichi Ishiwatari, Yuki Yasuda, Taro Miyazaki, and Jun Goto. 2020. Relation-aware Graph Attention Networks with Relational Position Encodings for Emotion Recognition in Conversations. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP.* 7360–7370. https://doi.org/10.18653/v1/2020.emnlp-main.597

[18] Abhinav Joshi, Ashwani Bhat, Ayush Jain, Atin Vikram Singh, and Ashutosh Modi. 2022. COGMEN: COntextualized GNN based Multimodal Emotion recognitioN. In

*Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL.* 4148–4164. https://doi.org/10.18653/v1/2022.naacl-main.306

[19] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR.* http://arxiv.org/abs/1412.6980

[20] Robert W Levenson. 2011. Basic emotion questions. *Emotion review* 3, 4 (2011), 379–386.

[21] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL.* 7871–7880. https://doi.org/10.18653/v1/2020.acl-main.703

[22] Shimin Li, Hang Yan, and Xipeng Qiu. 2022. Contrast and Generation Make BART a Good Dialogue Emotion Recognizer. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022.* 11002–11010. https://ojs.aaai.org/index.php/AAAI/article/view/21348

[23] Zaijing Li, Fengxiao Tang, Ming Zhao, and Yusen Zhu. 2022. EmoCaps: Emotion Capsule based Model for Conversational Emotion Recognition. In *Findings of the Association for Computational Linguistics: ACL.* 1610–1618. https://doi.org/10.18653/v1/2022.findings-acl.126

[24] Zheng Lian, Lan Chen, Licai Sun, Bin Liu, and Jianhua Tao. 2022. GCNet: Graph Completion Network for Incomplete Multimodal Learning in Conversation. *CoRR* abs/2203.02177 (2022). https://doi.org/10.48550/arXiv.2203.02177

[25] Yunze Liu, Qingnan Fan, Shanghang Zhang, Hao Dong, Thomas A. Funkhouser, and Li Yi. 2021. Contrastive Multimodal Fusion with TupleInfoNCE. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV.* 734–743. https://doi.org/10.1109/ICCV48922.2021.00079

[26] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR* abs/1907.11692 (2019). arXiv:1907.11692 http://arxiv.org/abs/1907.11692

[27] Hui Ma, Jian Wang, Hongfei Lin, Xuejun Pan, Yijia Zhang, and Zhihao Yang. 2022. A multi-view network for real-time emotion recognition in conversations. *Knowl. Based Syst.* 236 (2022), 107751. https://doi.org/10.1016/j.knosys.2021.107751

[28] Yukun Ma, Khanh Linh Nguyen, Frank Z. Xing, and Erik Cambria. 2020. A survey on empathetic dialogue systems. *Inf. Fusion* 64 (2020), 50–70. https://doi.org/10.1016/j.inffus.2020.06.011

[29] Navonil Majumder, Soujanya Poria, Devamanyu Hazarika, Rada Mihalcea, Alexander F. Gelbukh, and Erik Cambria. 2019. DialogueRNN: An Attentive RNN for Emotion Detection in Conversations. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019.* 6818–6825. https://doi.org/10.1609/aaai.

v33i01.33016818

[30] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, Volume 1 (Long Papers)*. 2227–2237. https://doi.org/10.18653/v1/n18-1202

[31] Soujanya Poria, Erik Cambria, Devamanyu Hazarika, Navonil Majumder, Amir Zadeh, and Louis-Philippe Morency. 2017. Context-Dependent Sentiment Analysis in User-Generated Videos. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL, Volume 1: Long Papers*. 873–883. https://doi.org/10.18653/v1/P17-1081

[32] Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2019. MELD: A Multimodal Multi-Party Dataset for Emotion Recognition in Conversations. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL, Volume 1: Long Papers*. 527–536. https://doi.org/10.18653/v1/p19-1050

[33] Weizhou Shen, Junqing Chen, Xiaojun Quan, and Zhixian Xie. 2021. DialogXL: All-in-One XLNet for Multi-Party Conversation Emotion Recognition. In *Thirty-Fifth AAAI Conference on Artificial Intelligence,*. 13789–13797. https://ojs.aaai.org/index.php/AAAI/article/view/17625

[34] Weizhou Shen, Siyue Wu, Yunyi Yang, and Xiaojun Quan. 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, ACL/IJCNLP 2021, (Volume 1: Long Papers)*. 1551–1560. https://doi.org/10.18653/v1/2021.acl-long.123

[35] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2019. Multimodal Transformer for Unaligned Multimodal Language Sequences. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL, Volume 1: Long Papers*. 6558–6569. https://doi.org/10.18653/v1/p19-1656

[36] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation Learning with Contrastive Predictive Coding. *CoRR* abs/1807.03748 (2018). arXiv:1807.03748 http://arxiv.org/abs/1807.03748

[37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems,NeurIPS*. 5998–6008. https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

[38] Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang, Wei Wu, and Weiran Xu. 2021. ConSERT: A Contrastive Framework for Self-Supervised Sentence Representation Transfer. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP, (Volume 1: Long Papers)*. 5065–5075. https://doi.org/10.18653/v1/2021.acl-long.393

[39] Jianing Yang, Yongxin Wang, Ruitao Yi, Yuying Zhu, Azaan Rehman, Amir Zadeh, Soujanya Poria, and Louis-Philippe Morency. 2021. MTAG: Modal-Temporal Attention Graph for Unaligned Human Multimodal Language Sequences. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*. 1009–1021. https://doi.org/10.18653/v1/2021.naacl-main.79

[40] Dong Zhang, Weisheng Zhang, Shoushan Li, Qiaoming Zhu, and Guodong Zhou. 2020. Modeling both Intra- and Inter-modal Influence for Real-Time Emotion Detection in Conversations. In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*. 503–511. https://doi.org/10.1145/3394171.3413949

[41] Lixing Zhu, Gabriele Pergola, Lin Gui, Deyu Zhou, and Yulan He. 2021. Topic-Driven and Knowledge-Aware Transformer for Dialogue Emotion Detection. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP*. 1571–1582. https://doi.org/10.18653/v1/2021.acl-long.125

[42] Shihao Zou, Xianying Huang, XuDong Shen, and Hankai Liu. 2022. Improving multimodal fusion with Main Modal Transformer for emotion recognition in conversation. *Knowl. Based Syst.* 258 (2022), 109978. https://doi.org/10.1016/j.knosys.2022.109978