

M³HG: Multimodal, Multi-scale, and Multi-type Node Heterogeneous Graph for Emotion Cause Triplet Extraction in Conversations

Qiao Liang Ying Shen* Tiantian Chen Lin Zhang

Tongji University, Shanghai, China

{2333091, yingshen, 2111287, cslinzhang}@tongji.edu.cn[†]

Abstract

Emotion Cause Triplet Extraction in Multimodal Conversations (MECTEC) has recently gained significant attention in social media analysis, aiming to extract emotion utterances, cause utterances, and emotion categories simultaneously. However, the scarcity of related datasets, with only one published dataset featuring highly uniform dialogue scenarios, hinders model development in this field. To address this, we introduce **MECAD**, the first multimodal, multi-scenario MECTEC dataset, comprising 989 conversations from 56 TV series spanning a wide range of dialogue contexts. In addition, existing MECTEC methods fail to explicitly model emotional and causal contexts and neglect the fusion of semantic information at different levels, leading to performance degradation. In this paper, we propose M³HG, a novel model that explicitly captures emotional and causal contexts and effectively fuses contextual information at both inter- and intra-utterance levels via a multimodal heterogeneous graph. Extensive experiments demonstrate the effectiveness of M³HG compared with existing state-of-the-art methods. The codes and dataset are available at <https://github.com/redifinition/M3HG>.

1 Introduction

Emotion Cause Analysis in Conversations (ECAC) aims at identifying emotions and their causes in conversations, which is a crucial research field in natural language processing (Li et al., 2022b; Wang et al., 2023). However, most of ECAC research (Li et al., 2022b; Wang et al., 2023; Zheng et al., 2023; Chen et al., 2023) only focuses on the textual contexts, overlooking other modalities (Soleymani et al., 2017).

* Corresponding authors.

[†]This work was supported in part by the National Natural Science Foundation of China under Grant 62476202 and 62272343, in part by the Fundamental Research Funds for the Central Universities.

To address this limitation, Wang et al. (2022) proposed a new task called Multimodal Emotion Cause Triplet Extraction in Conversations (MECTEC). The task aims to simultaneously identify the emotion utterance, the corresponding cause utterances, and the emotion category (i.e., the *utter-cause-emotion* triplet) from a conversation containing three modalities: text, audio, and video. Figure 1 illustrates a multimodal conversation between a mother and daughter. In this example, there are six non-neutral utterances, and consequently, six utter-cause-emotion triplets are identified. MECTEC differs from ECAC in 1) multimodal contexts (i.e., text, audio, and video) resulting in more complex emotional expression, and 2) multi-scale semantic information from overall conversation and utterance features like intonation and facial expressions, which pose significant challenges.

Another major challenge in MECTEC is the scarcity of datasets. While numerous text-based datasets exist for ECAC, only one dataset, namely the ECF dataset (Wang et al., 2022), is specifically designed for MECTEC. However, the videos in ECF are all from the *Friends* TV series with restricted speakers and scenarios, hindering MECTEC model development. Therefore, in this work, a new multimodal, multi-scenario MECTEC dataset, namely **MECAD**, is constructed. To the best of our knowledge, it is the first of its kind and will greatly facilitate research in this field.

Constrained by the limited dataset, existing MECTEC models have various deficiencies. Wang et al. (2022) proposed a two-stage architecture that predicts emotion and cause utterances separately. However, this approach is computationally intensive and prone to error accumulation. Therefore, recent studies (Hu et al., 2024; Wang et al., 2023; Li et al., 2024a) propose one-stage architectures using graph neural networks or prompt engineering to extract utter-cause-emotion triplets. However, these methods do not explicitly extract specific contexts

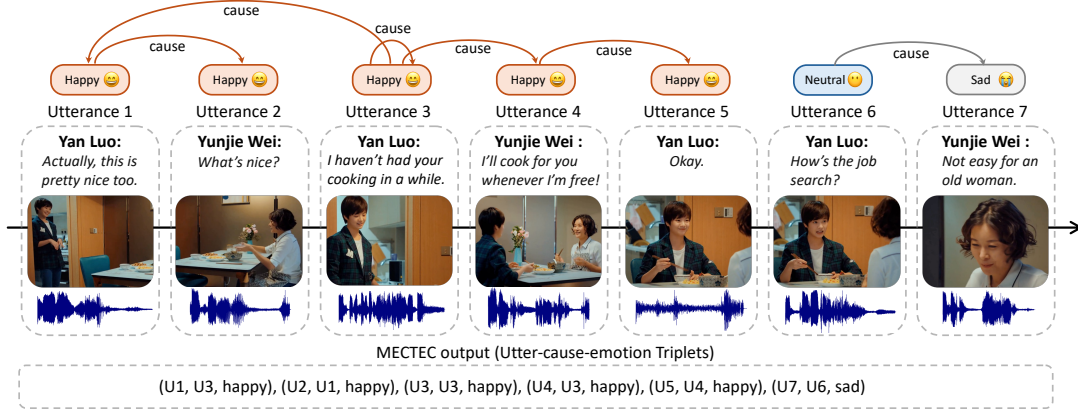


Figure 1: An example of the MECTEC task. Each utterance contains three different modalities - text, audio, and video. Arrows represent causal relationships that link the cause utterances to the corresponding emotion utterances. The dashed box at the bottom lists all the <utter-cause-emotion> triplets identified in this example.

related to emotions and their causes. According to emotion attribution theory (Weiner, 1985), the relationships of emotions and their causes are revealed by specific contexts, such as emotional words in texts, and intonations in audio and video conversations. For example, in Utterance 3 in Figure 1, a pleasant facial expression indicates happiness, while “haven’t had your cooking” and a happy tone reveal the cause. The example illustrates that emotions and their causes depend on contextual cues across multiple modalities, highlighting the necessity of explicitly modeling their specific contexts.

In addition, previous work (Wang et al., 2022; Hu et al., 2024; Wang et al., 2023; Li et al., 2024a; Wei et al., 2020) fail to effectively identify the **cause utterances occurring after emotion utterances**. For example, in Utterance 1 in Figure 1, the reason why Luo is happy cannot be obtained only from the historical context of Utterance 1. To find out the real cause of emotion in Utterance 1, the whole conversation should be scrutinized, which is overlooked by previous work.

Furthermore, existing models (Wang et al., 2022; Hu et al., 2024; Wang et al., 2023; Li et al., 2024a) fail to adequately extract semantic information at different scales. As shown in Figure 1, the semantic information that reveals the relationship of an utterance and its cause not only resides in inter-connections between utterances but also resides in the intra-content of each utterance. Therefore, it’s essential to comprehensively integrate semantic information in different scales during modality fusion.

To solve the aforementioned problems, we propose an MECTEC model based on the multimodal, multi-scale, and multi-type node heterogeneous

graph, named **M³HG**. M³HG accurately extracts emotion and cause-related contexts and fuses multimodal, multi-scale semantic information using multimodal heterogeneous graph attention network (HGAT) with multi-type nodes.

Our contributions can be summarized as follows:

- The first Chinese multi-scenario MECTEC dataset, **MECAD**, and an online sentiment data annotation toolkit are constructed. The dataset consists of 989 conversations with 10,519 utterances annotated with important information such as emotion labels, their causes, and types of emotional causes. It will greatly benefit the development of models in the MECTEC and related fields.
- An efficient MECTEC model, namely M³HG, is proposed to identify utter-cause-emotion triplets from multimodal conversations. It explicitly extracts specific emotion and cause-related contexts to find connections between emotions and causes. Besides, it fully integrates semantic information from inter and intra-utterance levels to enhance the model’s predictive ability.
- Extensive experiments are performed to verify the performance of our proposed model and other state-of-the-art models on MECAD and ECF datasets. Experimental results reveal that M³HG outperforms its counterparts, which demonstrates the effectiveness of our model.

2 Related Works

Emotion Cause Analysis in Conversations. Most existing studies on ECAC focus on Causal Emotion Entailment (CEE) and Emotion Cause Pair Extrac-

tion in Conversations (ECPEC). CEE aims to identify which cause utterances trigger the non-neutral emotions of the target utterances. Since CEE assumes emotion utterances are given, most related work (Poria et al., 2021; Li et al., 2022a; Zhang et al., 2022; Gu et al., 2023) viewed CEE as an utterance classification problem. However, because emotions of utterances are often unknown in real-world conversations, Li et al. (2022b) proposed the ECPEC task which additionally predicts emotions for the target utterances. Subsequent work (Wang et al., 2023; Zhao et al., 2023) has incorporated commonsense knowledge into GATs to improve the model’s semantic understanding of emotions and causes, achieving better performance. Besides, some methods (Ding et al., 2020a,b; Wei et al., 2020; Zheng et al., 2022) from models in the Emotion cause Pair Extraction (ECPE) field are also adapted for the ECPEC task.

Multimodal Emotion Cause Triplet Extraction in Conversations. In recent years, multimodal conversation scenarios on social media platforms have grown significantly, as more individuals share their lives and express emotions through live streaming and various online chats. To advance emotion cause analysis in multimodal conversation scenarios, Wang et al. (2022) introduced the MECTEC task and released the ECF dataset. However, few solutions have been proposed for this recently introduced task. Li et al. (2024a) incorporated emotion transition information into emotion-cause pair extraction using a novel labeling constraint, while Hu et al. (2024) fused semantic information across modalities via prompt engineering. These methods treat multimodal fusion and contextual information extraction for emotional causes as separate processes. Furthermore, they fail to effectively integrate semantic information across different scales, which significantly hampers the overall performance of models in the MECTEC task. To address these issues, we propose a model that fully integrates multi-scale semantic information from different modalities, preventing the loss of contextual information during fusion and improving triplet extraction accuracy.

Datasets for the ECAC Task. Table 1 summarizes popular datasets in ECAC. Poria et al. (2021) introduced the RECCON dataset for the ECAC task, and Li et al. (2022b) extended it by building the ConvECPE dataset. Given the multimodal nature of conversations, Wang et al. (2022) developed the ECF dataset for MECTEC. However, all scenes

Table 1: A summary of datasets for ECAC task. T, A, V stand for text, audio and video respectively.

Dataset	Modalities	Sources	# Instances
RECCON	T	Act and Daily	11,769
ConvECPE	T	Act	7,433
ECF	T,A,V	TV <i>Friends</i>	13,509
MECAD	T,A,V	56 TV series	10,516

in ECF are drawn from the *Friends*, limiting the diversity of conversation scenarios and contents.

3 Proposed MECAD Dataset

To facilitate the research in MECTEC and other related fields, we constructed a multi-scenario MECTEC dataset called MECAD. Compared with ECF (Wang et al., 2022), MECAD has more diverse conversation scenarios. In addition to labeling emotion categories and their causes for each utterance, we also categorized the types of emotion causes (e.g., *event*, *expression*) and the modality of annotation (i.e., *text*, *audio*, or *video*) to support future studies in multimodal emotion cause analysis.

We selected the publicly available M³ED (Zhao et al., 2022) dataset as our data source, which contains 990 segments from 56 Chinese TV series. However, M³ED dataset only contains conversation scripts, audios, and screenshots, lacking corresponding videos. Therefore, we endeavored to collect the corresponding video segments based on the conversation timestamps provided by M³ED. We concatenated sentences to form 989 multimodal conversations with 10,516 full utterances.

We invited 10 Chinese graduate students majored in Psychology to annotate the corresponding cause utterances, the types of emotion causes and the modal cues of annotations in the conversations. To obtain high-quality annotations, we designed detailed guidelines based on previous studies (Dirven, 1997; Steptoe and Brydon, 2009), trained the volunteers, and tested them with annotation cases. Only those passing the test participated in the final annotation process. Each volunteer was paid \$50 for their annotations. Then, we randomly assigned three qualified annotators for each conversation. If divergence exists among annotations from different volunteers, the final annotation for the utterance is determined by majority voting. Two strategies were used to review and revise incorrect annotations: 1) Annotation consistency among the three annotators for each TV series is calculated. For series with low consistency, the annotators rechecked

and revised their labels as needed. 2) If disagreements remained, a fourth annotator was invited to relabel the utterances and make the final decision.

To enhance annotation efficiency and accuracy, we developed an online multimodal conversation emotion cause annotation tool. The interface of the annotation tool is shown in Figure 3 in Appendix B. This tool is highly reusable and user-friendly, making it ideal for related research in the future.

We use Cohen’s Kappa (Cohen, 1960) to assess pairwise agreement and Fleiss’s Kappa (McHugh, 2012) for overall consistency among annotators. The Cohen’s Kappa results are in Appendix A, and the Fleiss’s Kappa score of 0.6932 exceeds the threshold of 0.61 (Landis, 1977), confirming the statistical reliability of our annotations.

The dataset statistics and detailed analysis of MECAD are presented in Figure 4 in Appendix A. MECAD provides solid support for assessing the performance and generalization capabilities of MECTEC models in broader scenarios.

4 Framework of Proposed M³HG

4.1 Task Definition

Given a conversation $C = \{(S_i, U_i)\}_{1 \leq i \leq n}$, where S_i denotes the speaker of the i -th utterance U_i , n denotes the length of the conversation C , $U_i = \{U_i^t, U_i^a, U_i^v\}$, and t, a, v are the text, audio and video modality, respectively. The goal of MECTEC is to identify all the utter-cause-emotion triplets from the conversation C :

$$\mathcal{P} = \{(U_j^e, U_j^c, y_j^e)\}, \quad (1)$$

where U_j^e is the j -th utterance with emotion y_j^e , U_j^c is the corresponding cause utterance, and $y_j^e \in \{\text{Anger, Disgust, Fear, Joy, Sadness, Surprise}\}$ (Ekman, 1992).

4.2 Model Overview

M³HG is an end-to-end (E2E) MECTEC model, as illustrated in Figure 2. It consists of four key components: *unimodal feature extraction*, *graph construction*, *multi-scale semantic fusion*, and *emotion-cause classification*.

In *unimodal feature extraction*, M³HG extracts local contextual representations for each utterance using modality-specific feature extractors and unimodal encoders. In *graph construction*, M³HG constructs a conversation interaction graph using these feature representations to explicitly model

the emotion and cause-related contexts. In *multi-scale semantic fusion*, M³HG combines semantic information at different scales within the conversation interaction graph to produce a comprehensive feature representation of both emotion and cause contexts. In *emotion-cause classification*, emotion and cause contextual representations are concatenated and used to extract utter-cause-emotion triplets with position embedding.

4.3 Unimodal Feature Extraction

First, we utilize SA-RoBERTa (Gu et al., 2020), Wav2Vec2 (Baevski et al., 2020), and DenseNet (Huang et al., 2017) to extract three feature representations E^t , E^a , and E^v , from text, audio, and video, respectively, where $E^t \in \mathbb{R}^{n \times d_t}$, $E^a \in \mathbb{R}^{n \times d_a}$, and $E^v \in \mathbb{R}^{n \times d_v}$, and d_t, d_a, d_v represent dimensions of the hidden layer representations of the three modalities. The extraction process is described in Appendix C.1.

Then, we encode each feature representation within an unimodal local context. For text, we apply multi-head self-attention (Vaswani, 2017) to E^t to capture local contextual information, resulting in H^t . For E^a and E^v , we use a GRU-based network (Li et al., 2024b) to extract local context by leveraging the RNN structure’s capability to handle temporal features, which is expressed as:

$$\begin{aligned} E'^m &= LN(E^m + GRU(E^m)), \\ H^m &= LN(E^m + E'^m + FFN(E'^m)), \end{aligned} \quad (2)$$

where $H^m \in \mathbb{R}^{n \times d_m}$, $m \in \{a, v\}$, LN denotes layer normalization, and FFN denotes a feedforward neural network.

After encoding the local context for each modality, we obtain the sequence representations H^t , H^a , H^v for text, audio, and video. We then apply three linear layers to map H^t , H^a , H^v to H'^t , H'^a , H'^v with the same dimension d_h .

4.4 Graph Construction

To enable M³HG to fuse multi-scale semantic information across modalities, we construct a heterogeneous graph that represents both inter- and intra-utterance connections, as well as cross-modal interactions. The structure of this heterogeneous graph can be denoted by $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{R})$, where \mathcal{V} is the node set consisting of all graph nodes v_i , \mathcal{R} is the relation set consisting of all relations r_{ij} between any two nodes v_i and v_j , and \mathcal{E} is the edge set consisting of all edges represented as (v_i, r_{ij}, v_j) .

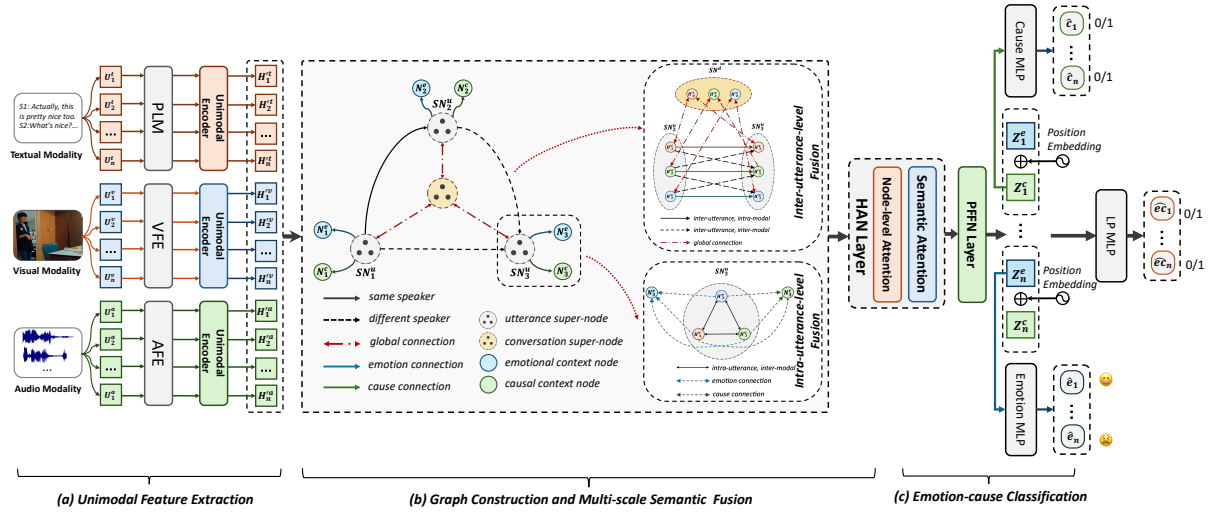


Figure 2: The framework of proposed M³HG. It consists of three main components: unimodal feature extraction, graph construction and multi-scale semantic fusion, and emotion-cause classification.

Nodes. To explicitly model emotion and cause-related contexts in conversations, we model them as *emotional context nodes* N^e and *causal context nodes* N^c , respectively. To enable \mathcal{G} to accurately perceive the conversation information, we model the whole conversation as a *conversation node*. Each utterance is represented by an *utterance node*. Both the utterance node and conversation node are designed as Super-Nodes containing these modalities, denoted as SN^u and SN^d , since they contain three modal features. Therefore, \mathcal{G} contains four types of nodes: N^e , N^c , SN^u and SN^d .

N^e and N^c are first initialized with textual sequence representations H^t , then updated with contextual information from the other two modalities, which is described in Section 4.5. Each utterance Super-Node $SN^u = \{N^t, N^a, N^v\}$ is initialized using H^t, H^a, H^v . The conversation node $SN^d = \{N_d^t, N_d^a, N_d^v\}$ is initialized by averaging H^t, H^a, H^v to capture global information.

Edges and Relations. There are five types of Super-Edges connecting the aforementioned Super-Nodes: *same speaker* (r_{ss}), *different speaker* (r_{ds}), *global connection* (r_{gc}), *emotion connection* (r_{ec}) and *cause connection* (r_{cc}). The *same speaker* edge connects the utterance Super-Nodes SN^u from the same speaker. Inspired by the work of Shen et al. (2021), we define the local context as K preceding utterances from the same speaker of SN^u , where K is a hyper-parameter. The *different speaker* edge connects the utterance Super-Nodes within the local context from different speakers to SN^u . The

bidirectional *global connection* edge connects all the utterance Super-Nodes SN^u s with the conversation Super-Node SN^d , facilitating the propagation of global contextual information. The *emotion connection* edge and the *cause connection* edge connect SN^u with its corresponding emotional context node N^e and causal context node N^c , respectively. They explicitly capture the emotion and cause context specific to each utterance.

M³HG is the first MECTEC model capable of handling situations where **cause utterances appear after emotion utterances**, as each utterance is linked through the *global connection* node. The detailed experiments in Appendix E.2 further validate this capability. The pseudo-code of graph construction and a constructed graph for the conversation in Figure 1 are provided in Appendix C.2 and Appendix C.3, respectively. The graph construction process of M³HG can be expressed as:

$$\begin{aligned}
 \mathcal{G} &= (\mathcal{V}, \mathcal{E}, \mathcal{R}), \\
 \mathcal{V} &= \{SN_i^u, N_i^e, N_i^c, SN^d\}_{1 \leq i \leq n}, \\
 SN_i^u &= \{N_i^t, N_i^a, N_i^v\}, \\
 SN^d &= \{N_d^t, N_d^a, N_d^v\}, \\
 \mathcal{R} &= \{r_{ss}, r_{ds}, r_{gc}, r_{ec}, r_{cc}\}, \\
 \mathcal{E} &= \{(v_i, r_{ij}, v_j)\}, v_i, v_j \in \mathcal{V}, r \in \mathcal{R},
 \end{aligned} \tag{3}$$

where superscripts u, e, c, d denote node types, and m denotes three modalities. Based on the constructed graph \mathcal{G} , the emotion and cause contexts are effectively modeled.

4.5 Multi-scale Semantic Information Fusion

Based on graph \mathcal{G} , we designed a comprehensive approach to integrate semantic information across different modalities and scales. This mechanism is implemented in two levels: **intra-utterance fusion** which captures emotion and cause-related contexts within utterances, and **inter-utterance fusion** which propagates semantic information among utterances and conversation-level contexts. Both levels leverage HGAT (Wang et al., 2019) to propagate and fuse semantic information through various meta-paths (Wang et al., 2019) within \mathcal{G} . This ensures thorough updates to node features by integrating multi-scale semantic information.

The meta-paths in \mathcal{G} are defined as:

$$\begin{aligned} \Phi &= \{\phi(v_i, r_{ij}, v_j)\}, v_i, v_j \in \mathcal{V}, \\ \phi(v_i, r_{ij}, v_j) &= v_i \xleftrightarrow{r_{ij}} v_j, r_{ij} \in \mathcal{R}, \end{aligned} \quad (4)$$

where $\phi(v_i, r_{ij}, v_j)$ represents all paths that connect node v_i to node v_j via edge type r_{ij} .

Intra-utterance-level Fusion. As shown in Figure 2, for each utterance Super-Node SN^u , we perform intra-utterance-level fusion by integrating semantic information within the utterance. We define the meta-path Φ_{intra} for intra-utterance-level semantic fusion for SN_n^u as:

$$\begin{aligned} \Phi_{intra} &= \{\phi(N^{m_1}, N^{m_2}, r_{m_1, m_2})\} \\ &\cup \{\phi(N^m, N^e, r_{m, e})\} \\ &\cup \{\phi(N^m, N^c, r_{m, c})\}, \end{aligned} \quad (5)$$

where $m_1, m_2, m \in \{t, a, v\}$, N^m represents the nodes of modality m within the SN^u , and r_{m_1, m_2} denotes the edges connecting N^{m_1} and N^{m_2} . $r_{m, e}$ denotes edges connecting nodes N^m to the emotional context nodes N^e , facilitating the aggregation of emotional contexts conveyed by different modalities within the utterance. Similarly, $r_{m, c}$ represents the edges that connect N^m to the causal context nodes N^c , enabling the aggregation of causal contexts. Φ_{intra} effectively models the process of semantic information fusion in a single utterance.

Next, we incorporate node-level attention into Φ_{intra} . For each meta-path in Φ_{intra} and nodes $v_i \in \{N^m, N^e, N^c\}$, the importance of its neighbors \mathcal{N}_i in Φ_{intra} is computed as:

$$\alpha_{ij}^\phi = \frac{\exp(\sigma(\mathbf{a}_\phi^T \cdot [\mathbf{H}'_i \parallel \mathbf{H}'_j]))}{\sum_{k \in \mathcal{N}_i^\phi} \exp(\sigma(\mathbf{a}_\phi^T \cdot [\mathbf{H}'_i \parallel \mathbf{H}'_k]))}, \quad \phi \in \Phi_{intra}, \quad (6)$$

where σ denotes the activation function, and \mathbf{a}_ϕ is the node-level attention vector of meta-path ϕ . The node representation of v_i based on meta-path ϕ is obtained by:

$$\mathbf{Z}_i = \sigma\left(\sum_{j \in \mathcal{N}_i^\phi} \alpha_{ij}^\phi \cdot \mathbf{H}'_j\right). \quad (7)$$

This process yields the contextual features $\mathbf{Z}_i \in \mathbb{R}^{1 \times d_h}$ for nodes v_i under the intra-utterance-level meta-paths Φ_{intra} .

Inter-utterance-level Fusion. As illustrated in Figure 2, for any two utterance Super-Nodes SN_i^u and SN_j^u in \mathcal{G} , along with the conversation Super-Node SN^d , we perform inter-utterance-level fusion by connecting SN_i^u and SN_j^u to SN^d , thereby integrating contextual information across utterances. We define meta-paths Φ_{inter} for inter-utterance-level fusion between SN^u and SN^d :

$$\begin{aligned} \Phi_{inter} &= \{\phi(N_i^{m_1}, N_j^{m_2}, r_{m_1, m_2})\} \\ &\cup \{\phi(N_i^m, N_d^m, r_{d, m})\} \\ &\cup \{\phi(N_j^m, N_d^m, r_{d, m})\}, \end{aligned} \quad (8)$$

where $m_1, m_2, m \in \{t, a, v\}$, N_i^m and N_j^m represent the nodes of modality m inside SN_i^u and SN_j^u , respectively, r_{m_1, m_2} denotes the edges connecting $N_i^{m_1}$ and $N_j^{m_2}$, and $r_{d, m}$ represents the edges connecting SN^u s to SN^d in modality m .

The utterance information from each modality can be passed to SN^d through Φ_{inter} , which accomplishes inter-utterance-level fusion between utterances. As a result, SN^d comprehensively integrates information across all three modalities. The meta-path set Φ_{inter} models multimodal connections between utterances, enabling conversation information aggregated in \mathcal{G} .

Similar to Eq. 6 and Eq. 7, the contextual representations of SN^u and SN^d are obtained under the meta-path Φ_{inter} by the node-level attention block.

After performing multi-scale semantic fusion with Φ_{intra} and Φ_{inter} , we apply the semantic attention mechanism (Wang et al., 2019) to each node embedding \mathbf{Z}_i , integrating multi-scale semantic information from all three modalities. Following (Chen et al., 2023), each fusion iteration is followed by a position-wise feed-forward network (PFFN) layer, which updates node features through a non-linear transformation. The emotional context node representation \mathbf{Z}_i^e and the causal context node feature representation \mathbf{Z}_i^c can be obtained at the end of iterations of the multi-scale semantic fusion and PFFN layers.

4.6 Emotion-cause classification

For each utterance U_i , its Z_i^e and Z_i^c are fed into the emotion-specific Multi-Layer Perceptron (Emotion MLP) and the cause-specific Multi-Layer Perceptron (Cause MLP) to predict its emotion category \hat{y}_i^e and the cause indicator \hat{y}_i^c which indicates whether U_i can be a cause utterance. For each utterance pair U_i and U_j , we compute a relative position encoding RPE_{ij} to capture the positional relationship between U_i and U_j . We utilize the RBF kernel function (Wei et al., 2020) to compute RPE_{ij} , which captures the relative positional relationships between utterances through a nonlinear relation. Z_j^e , Z_i^c and RPE_{ij} are then concatenated and fed into a new MLP to determine whether U_i is the cause utterance of U_j :

$$\hat{y}_{ij}^{ec} = \sigma(MLP(Z_j^e || Z_i^c || RPE_{ij})). \quad (9)$$

\hat{y}_{ij}^{ec} represent the binary classification logits indicating whether U_i is the cause of U_j . Based on \hat{y}_{ij}^{ec} , we can determine whether U_j , U_i and \hat{y}_j^e can form a true utter-cause-emotion triplet.

4.7 Training

We use Focal loss (Ross and Dollár, 2017) to cope with category imbalance in emotion-cause classification. Specifically, the loss of both emotion prediction and cause utterance prediction and the emotion-cause pair prediction, can be expressed as:

$$\mathcal{L}^\beta = -\frac{1}{N^\beta} \sum_{i=1}^{N^\beta} \alpha^\beta (1 - \hat{y}_i^\beta)^\gamma \log(\hat{y}_i^\beta), \beta \in \{e, c, ec\} \quad (10)$$

where β represents the task type, N^β denotes the corresponding sample number of β , α^β is the category balancing factor, and γ denotes the Focal loss modulation parameter. These three training losses are optimized jointly during the training process.

5 EXPERIMENTS

5.1 Experimental Settings

We conduct extensive experiments on two MECTEC benchmark datasets, i.e., **ECF** (Wang et al., 2022) and **MECAD**, which both contain data of three modalities: text, audio, and video. Similar to (Wang et al., 2022), we evaluate the model’s overall performance using the F1 score. The F1 score is computed for utter-cause-emotion triplets within each emotion category separately. Then the weighted average F1 score is calculated across all

six emotion categories which is referred to 6 Avg. In addition, as in (Wang et al., 2023), considering the data imbalance among different emotion categories, we also report the weighted average F1 scores for the four main emotion categories except *Disgust* and *Fear*, which is referred to 4 Avg. The implementation details of the experiment are given in Appendix D.

5.2 Baselines

Due to the limited research on the MECTEC task, representative approaches in related fields of Emotion Cause Pair Extraction (ECPE) and Emotion Cause Pair Extraction in Conversations (ECPEC) are considered. The ECPE and ECPEC tasks aim to extract emotion-cause pairs from plain texts and conversations, respectively.

We compare our model with seven baselines: 1) **MC-ECPE-2steps** (Wang et al., 2022) is a two-step MECTEC architecture, which first extracts emotion utterances and cause utterances separately, and then performs pairing and filtering to identify emotion-cause pairs. 2) **HiLo** (Li et al., 2024a) is one of the SOTA approaches for the MECTEC task, which fully utilizes conversation information through a labeling constraint mechanism. 3) **ECPE-2D** (Ding et al., 2020a) is an E2E framework for ECPE that uses 2D-Transformer to model the interactions of emotion-cause pairs. 4) **RankCP** (Wei et al., 2020) is a GAT-based approach for ECPE to extract emotion-cause pairs by ranking. 5) **UECA-Prompt** (Zheng et al., 2022) is one of the SOTA methods for ECPE, which decomposes the task into multiple objectives and converts them into sub-prompts. 6) **SHARK** (Wang et al., 2023) is the SOTA method for ECPEC that incorporates commonsense into GATs to improve the model’s semantic understanding of emotions and causes. 7) **GPT-4o** is one of the most powerful large language models (LLMs) for open-domain conversations. Details of prompts are provided in Appendix F.

5.3 Experimental Results

Table 2 shows the experimental results of M³HG and seven baseline models evaluated on the ECF dataset and the MECAD dataset. Our model demonstrates an excellent performance both on the ECF dataset and the MECAD dataset.

Results on the ECF dataset. First of all, as shown in Table 2, among all the baseline models, the E2E approaches such as SHARK and HiLo deliver the

Dataset		Method	Modality	Anger	Disgust	Fear	Joy	Sadness	Surprise	6 Avg.	4 Avg.
ECF	Pipline	MC-ECPE-2steps [△]	T, A, V	24.39	0.00	0.71	38.84	21.60	40.24	29.32	31.92
	E2E	ECPE-2D [△]	T	25.13	0.00	0.00	41.25	21.62	43.24	30.80	33.55
		RankCP	T	28.29	12.03	3.52	38.69	22.17	37.67	30.58	32.48
		UECA-Prompt [△]	T	27.37	12.85	7.91	37.96	22.51	39.53	30.75	32.49
		SHARK [*]	T	28.65	10.42	5.33	40.41	25.35	40.45	32.24	34.33
		HiLo [*]	T, A, V	-	-	-	-	-	-	33.04	35.81
	LLMs	GPT-4o (5-shots)	T	28.49	17.76	12.35	31.11	27.27	33.89	29.13	30.30
	M ³ HG (ours)	T	34.47	18.17	12.72	43.28	32.22	45.82	37.46	39.95	
		T, A	35.53	18.71	17.07	47.73	30.97	46.72	39.10	40.97	
		T, V	34.05	18.18	19.57	46.23	32.10	48.50	38.90	40.72	
T, A, V		36.08	23.33	9.88	49.03	32.41	47.46	40.07	41.96		
MECAD	Pipeline	MC-ECPE-2steps	T, A, V	28.43	0.00	0.23	22.45	27.67	45.14	22.01	24.83
	E2E	ECPE-2D	T	28.12	0.00	0.56	24.30	28.01	35.87	25.32	28.54
		RankCP	T	29.79	12.50	3.06	21.79	29.31	32.36	26.29	28.32
		UECA-Prompt	T	28.54	12.12	5.32	20.84	29.67	34.17	25.91	27.87
		SHARK	T	30.22	10.16	4.10	25.84	30.21	34.59	27.58	29.99
		HiLo [*]	T, A, V	-	-	-	-	-	-	-	-
	LLMs	GPT-4o (5-shots)	T	36.65	20.08	8.45	24.52	17.89	39.77	27.16	28.42
	M ³ HG (ours)	T	35.85	18.05	15.38	25.95	29.13	42.11	30.81	32.55	
		T, A	37.29	21.03	15.89	27.15	30.34	42.78	32.16	33.73	
		T, V	36.91	20.48	16.91	25.47	30.96	43.14	31.95	33.52	
T, A, V		38.34	21.89	8.79	28.10	31.17	43.29	32.82	34.59		

Table 2: Performance comparison of different methods on the MECATEC task. [△] denotes the results are from (Wang et al., 2023). * denotes the results are from the original paper (Wang et al., 2023; Li et al., 2024a). The best results and the second best results are in bold and underlined, respectively. Since HiLo (Li et al., 2024a) is not publicly available, we only report the results of HiLo on the ECF dataset.

best performance, indicating that the E2E framework is more effective compared to the two-step pipeline frameworks. In contrast, M³HG adopting three modalities outperforms the SOTA E2E model HiLo, with 21.28% and 17.17% improvement in 6 Avg and 4 Avg scores, respectively. We attribute this improvement to M³HG’s ability to effectively extract semantic information at inter-utterance and intra-utterance levels, which enables the model to accurately pair emotion utterances and cause utterances. Specifically, in two challenging emotion categories which have limited training samples, i.e. *Disgust* and *Fear*, M³HG also exhibits high performances. For example, compared to GPT-4o, which achieved the second highest F1 scores in the *Disgust* and *Fear* categories, M³HG shows improvements of 31.36% and 58.46%, respectively.

When only incorporating the text modality, the 6 Avg and the 4 Avg scores of M³HG are 37.46 and 39.95. When incorporating audio and video with the text modality separately, the performance of M³HG is improved to 39.10, 38.90 of 6 Avg scores and 40.97, 40.72 of 4 Avg scores. When incorporating all three modalities, M³HG achieves the highest performance with 40.07 of 6 Avg scores and 41.96 of 4 Avg scores. Meanwhile, it can be observed that M³HG outperforms all the baseline models even when only using the text modality, demonstrating its superiority on the ECF dataset.

Results on the MECAD dataset. As shown in

Table 2, M³HG also achieves the highest results on the MECAD dataset. Compared to the second best model SHARK, M³HG adopting three modalities achieves the improvement of 19% on the 6 Avg scores and 15.34% on the 4 Avg score. Furthermore, despite GPT-4o’s superior semantic comprehension abilities, its performance on the MECAD dataset remains suboptimal, with its 6 Avg score and 4 Avg score of 27.16 and 28.42. Therefore, the few-shot-based LLM approach still struggles to effectively handle the MECATEC task. As shown in Table 2, M³HG exhibits a universal highest performance on the MECAD dataset, demonstrating the superiority and robustness of M³HG when dealing with multi-conversation scenarios.

More detailed experimental results and the ablation study on M³HG are presented in Appendix E.

6 Conclusion

In this work, we propose the first multimodal and multi-scenario Chinese emotion-cause analysis dataset, MECAD, for MECATEC and related emotion cause analysis tasks. Compared to ECF, the only existing dataset for multimodal emotion-cause analysis, MECAD offers more diverse conversation scenarios. It helps to enhance the generalizability and applicability of MECATEC models in complex social media environments. Moreover, MECAD is a valuable resource for cross-cultural emotion analysis and recognition. Furthermore, we propose

a generalized MECTEC framework named M³HG, which deeply extracts emotional and causal contexts, while effectively integrating semantic information across multiple granular levels. Extensive experiments on the ECF dataset and the MECAD dataset demonstrate the superiority of our method compared to the existing state-of-the-art methods.

Limitations

There are also some potential limitations in this work. First, the process of emotional and causal context extraction does not integrate external knowledge, which limits the model’s accuracy for emotion prediction and cause prediction. In the future, we plan to integrate external knowledge into our model and leverage the advanced semantic extraction capabilities of current LLM technology to facilitate deeper and more precise emotion cause analysis. Second, M³HG cannot handle excessively long conversations, as its input length is constrained by the language model used. Furthermore, M³HG may suffer from error propagation in the multimodal fusion process when emotion labels have uneven information across modalities. This imbalance can lead to inaccurate predictions, especially when modalities conflict. This challenge is common in current multimodal models for emotion-cause analysis and suggests an area for future improvement.

Ethical Considerations

We did not use real-world conversations in our data collection because such conversations may violate the privacy of the speaker. The effect of recruiting actors to play the roles is the same as in the TV series, but the scenes are not as diverse as in the TV series. Therefore, we use TV series as the data source. To further protect privacy, all data annotations were anonymized and de-identified, ensuring that our data collection adheres to ethical standards.

References

- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
- Tiantian Chen, Ying Shen, Xuri Chen, Lin Zhang, and Shengjie Zhao. 2023. Mpeg: A multi-perspective enhanced graph attention network for causal emotion

entailment in conversations. *IEEE Transactions on Affective Computing*.

Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.

Zixiang Ding, Rui Xia, and Jianfei Yu. 2020a. 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*, pages 3161–3170.

Zixiang Ding, Rui Xia, and Jianfei Yu. 2020b. End-to-end emotion-cause pair extraction based on sliding window multi-label learning. In *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*, pages 3574–3583.

René Dirven. 1997. Emotions as cause and the cause of emotions. *The language of emotions: Conceptualization, expression, and theoretical foundation*, pages 55–83.

Paul Ekman. 1992. An argument for basic emotions. *Cognition & emotion*, 6(3-4):169–200.

Jia-Chen Gu, Tianda Li, Quan Liu, Zhen-Hua Ling, Zhiming Su, Si Wei, and Xiaodan Zhu. 2020. Speaker-aware bert for multi-turn response selection in retrieval-based chatbots. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 2041–2044.

Xiaojie Gu, Renze Lou, Lin Sun, and Shangxin Li. 2023. Page: A position-aware graph-based model for emotion cause entailment in conversation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Guimin Hu, Zhihong Zhu, Daniel Hershcovich, Hasti Seifi, and Jiayuan Xie. 2024. Unimeec: Towards unified multimodal emotion recognition and emotion cause. *arXiv preprint arXiv:2404.00403*.

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*, pages 4700–4708.

JR Landis. 1977. The measurement of observer agreement for categorical data. *Biometrics*.

Bobo Li, Hao Fei, Fei Li, Tat-seng Chua, and Donghong Ji. 2024a. Multimodal emotion-cause pair extraction with holistic interaction and label constraint. *ACM Transactions on Multimedia Computing, Communications and Applications*.

Jiang Li, Xiaoping Wang, Yingjian Liu, and Zhigang Zeng. 2024b. Cfn-esa: A cross-modal fusion network with emotion-shift awareness for dialogue emotion recognition. *IEEE Transactions on Affective Computing*.

- Jiangnan Li, Fandong Meng, Zheng Lin, Rui Liu, Peng Fu, Yanan Cao, Weiping Wang, and Jie Zhou. 2022a. Neutral utterances are also causes: Enhancing conversational causal emotion entailment with social commonsense knowledge. *arXiv preprint arXiv:2205.00759*.
- Wei Li, Yang Li, Vlad Pandelea, Mengshi Ge, Luyao Zhu, and Erik Cambria. 2022b. Ecpec: Emotion-cause pair extraction in conversations. *IEEE Transactions on Affective Computing*, 14(3):1754–1765.
- I Loshchilov. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Mary L McHugh. 2012. Interrater reliability: the kappa statistic. *Biochemia medica*, 22(3):276–282.
- Soujanya Poria, Navonil Majumder, Devamanyu Hazarika, Deepanway Ghosal, Rishabh Bhardwaj, Samson Yu Bai Jian, Pengfei Hong, Romila Ghosh, Abhinaba Roy, Niyati Chhaya, et al. 2021. Recognizing emotion cause in conversations. *Cognitive Computation*, 13:1317–1332.
- T-YLPG Ross and GKHP Dollár. 2017. Focal loss for dense object detection. In *proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2980–2988.
- Weizhou Shen, Siyue Wu, Yunyi Yang, and Xiaojun Quan. 2021. Directed acyclic graph network for conversational emotion recognition. *arXiv preprint arXiv:2105.12907*.
- Mohammad Soleymani, David Garcia, Brendan Jou, Björn Schuller, Shih-Fu Chang, and Maja Pantic. 2017. A survey of multimodal sentiment analysis. *Image and Vision Computing*, 65:3–14.
- Andrew Steptoe and Lena Brydon. 2009. Emotional triggering of cardiac events. *Neuroscience & Biobehavioral Reviews*, 33(2):63–70.
- A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*.
- Fanfan Wang, Zixiang Ding, Rui Xia, Zhaoyu Li, and Jianfei Yu. 2022. Multimodal emotion-cause pair extraction in conversations. *IEEE Transactions on Affective Computing*, 14(3):1832–1844.
- Fanfan Wang, Jianfei Yu, and Rui Xia. 2023. Generative emotion cause triplet extraction in conversations with commonsense knowledge. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3952–3963.
- Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, and Philip S Yu. 2019. Heterogeneous graph attention network. In *The world wide web conference*, pages 2022–2032.
- Penghui Wei, Jiahao Zhao, and Wenji Mao. 2020. Effective inter-clause modeling for end-to-end emotion-cause pair extraction. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 3171–3181.
- Bernard Weiner. 1985. An attributional theory of achievement motivation and emotion. *Psychological review*, 92(4):548.
- Duzhen Zhang, Zhen Yang, Fandong Meng, Xiuyi Chen, and Jie Zhou. 2022. Tsam: A two-stream attention model for causal emotion entailment. *arXiv preprint arXiv:2203.00819*.
- Jinming Zhao, Tenggao Zhang, Jingwen Hu, Yuchen Liu, Qin Jin, Xinchao Wang, and Haizhou Li. 2022. M3ed: Multi-modal multi-scene multi-label emotional dialogue database. *arXiv preprint arXiv:2205.10237*.
- Weixiang Zhao, Yanyan Zhao, Zhuojun Li, and Bing Qin. 2023. Knowledge-bridged causal interaction network for causal emotion entailment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 14020–14028.
- Li Zheng, Donghong Ji, Fei Li, Hao Fei, Shengqiong Wu, Jingye Li, Bobo Li, and Chong Teng. 2023. Ecqed: emotion-cause quadruple extraction in dialogs. *arXiv preprint arXiv:2306.03969*.
- Xiaopeng Zheng, Zhiyue Liu, Zizhen Zhang, Zhaoyang Wang, and Jiahai Wang. 2022. Ueca-prompt: Universal prompt for emotion cause analysis. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 7031–7041.

A Dataset Statistics and Analysis of MECAD

To ensure the annotation quality of MECAD, we calculated Cohen’s kappa (Cohen, 1960) scores for every co-annotated data between two annotators, as shown in Figure 4. The Cohen’s kappa (Cohen, 1960) scores across all annotators are consistently around 0.6, indicating a good level of annotation consistency.

After the labeling was completed, we computed Cohen’s kappa scores separately for data that were not co-labeled between the two labelers, as shown in Figure 4. Table 3 lists some statistics of the MECAD dataset. The dataset contains a total of 989 conversations, 10,516 utterances, and 8,077 emotion cause pairs from 56 different TV series, which ensures the size and diversity of the dataset. Similar to M³ED (Zhao et al., 2022), we used TV-independent data segmentation to ensure the ability to validate model robustness as a benchmark dataset. The average number of utterances and the average length of an utterance of a conversation are similar in the training, validation, and test sets. At the same time, we can find that the average relative positions of the emotion cause pairs are all around

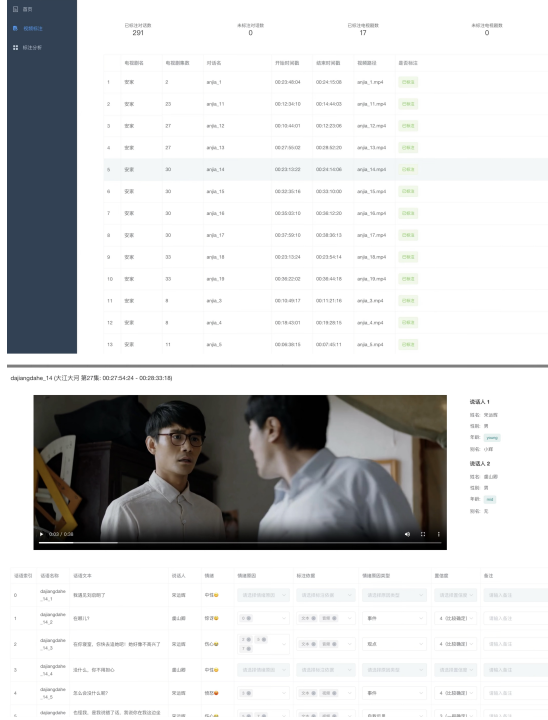


Figure 3: The interface of the developed online multi-modal conversation emotion cause annotation toolkit.

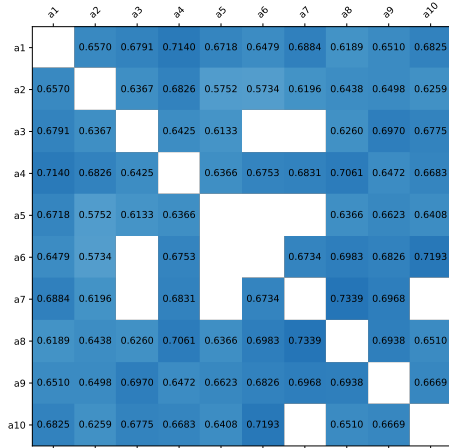


Figure 4: Schematic representation of Cohen's Kappa scores for the common labeled portion between every two annotators. A blank section indicates that there is no common annotation data between two annotators.

Table 3: MECAD statistics. *Rel pos of ec pairs* denotes the relative position between emotion utterances and cause utterances in emotion-cause pairs.

Statistic	Train	Val	Test	Total
# TV series	38	7	11	56
# conversations	684	126	179	989
# uttrs	7,516	1,168	1,832	10,516
# spkrs	421	87	118	626
Avg. uttrs/conversation	10.99	9.27	10.24	10.63
Avg. uttr length	18.30	18.80	18.15	18.33
Avg. rel pos of ec pairs	0.72	0.73	0.55	0.69
Max. rel pos of ec pairs	13	7	6	13
Min. rel pos of ec pairs	-14	-5	-9	-14
Emotion uttrs with cause ec pairs	4,526	743	1,062	6,331
	5,788	977	1,312	8,077

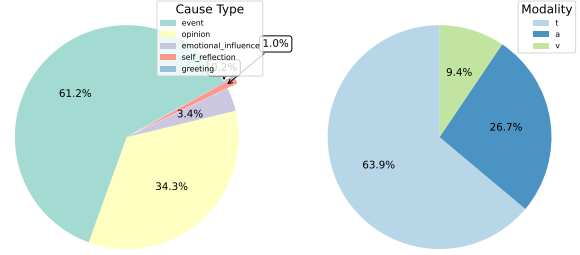


Figure 5: Percentage of five cause types in the MECAD dataset and percentage of modal basis for emotion cause inferences.

1, indicating that most of the emotions in the conversation are caused by the previous utterance.

We referred how the ECF (Wang et al., 2022) dataset categorizes the emotion causes and added a new category called *Self Reflection*, which differs from the remaining four categories by indicating that emotions may be triggered by an individual's introspection or self-reflection, such as recollections of past events or worries about the future. As shown in Figure 5, the event type is the cause type with the largest share, indicating that most of the emotions are caused by specific events in the conversation. Notably, 36.1% of the causes of emotion in our dataset are reflected in both audio and video modalities, which exemplifies the need for multimodal scene studies.

B The Annotation Toolkit of MECAD

To enhance annotation efficiency and accuracy, we developed an online multimodal conversation emotion cause annotation tool based on web technology¹. As illustrated in Figure 3, the toolkit's homepage presents a list of conversations assigned to the corresponding annotators, along with the progress of their annotations. The conversation an-

¹The annotation tool has been open-sourced at <https://github.com/redifinition/MECAD-MECTEC>

notation page displays speaker information, video segments, corresponding scripts, and configurable annotation items, enabling annotators to quickly and efficiently complete their annotations.

With flexible and modifiable web pages, researchers can utilize our annotation tools in dataset constructions for further multimodal sentiment analysis studies.

Algorithm 1 Super-Node-based Graph Construction for a Conversation

```

1: Input: the conversation  $\{S_1 : U_1, S_2 : U_2, \dots, S_N : U_N\}$ , speaker identity  $p(\cdot)$  satisfies  $p(U_i) = S_i$ , the direct context window  $K$ 
2: Output: Super-Node-based M3HG:  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{R})$ 
3:  $\mathcal{V} \leftarrow \{(SN_u^1, N_e^1, N_c^1), \dots, (SN_u^N, N_e^N, N_c^N), SN_d^1\}$ 
4:  $\mathcal{E} \leftarrow \emptyset$ 
5:  $\mathcal{R} \leftarrow \{r_{ss}, r_{ds}, r_{gc}, r_{ec}, r_{cc}\}$ 
6: for  $i \in \{2, 3, \dots, N\}$  do
7:    $c \leftarrow 0, w \leftarrow i - 1$ 
8:   while  $w > 0$  and  $c < K$  do
9:     if  $p(U_w) = p(U_i)$  then
10:       $\mathcal{E} \leftarrow \mathcal{E} \cup \{(SN_u^w, SN_u^i, r_{ss})\}$ 
11:       $c \leftarrow c + 1$ 
12:     else
13:       $\mathcal{E} \leftarrow \mathcal{E} \cup \{(SN_u^w, SN_u^i, r_{ds})\}$ 
14:     end if
15:      $w \leftarrow w - 1$ 
16:   end while
17: end for
18: for  $i \in \{1, 2, \dots, N\}$  do
19:    $\mathcal{E} \leftarrow \mathcal{E} \cup \{(SN_u^i, N_e^i, r_{ec})\}$ 
20:    $\mathcal{E} \leftarrow \mathcal{E} \cup \{(SN_u^i, N_c^i, r_{cc})\}$ 
21:    $\mathcal{E} \leftarrow \mathcal{E} \cup \{(SN_u^i, SN_d^i, r_{gc})\}$ 
22: end for
23: return  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{R})$ 

```

C Design Details of M³HG

C.1 Multimodal Feature Extracting

Text : We splice all the textual modal utterances and the corresponding speakers in the conversation and add a number of special tokens to get the textual modal input sequence: $X^t = \{< cls_token > S_1 : U_1^t, < sep_token >, \dots, < cls_token > S_n : U_n^t, < sep_token >\}$, where $< cls_token >$ and $< sep_token >$ denote the classification token and the separation token used in the pre-trained language model (PLM), respectively. To allow conversations that exceed the maximum input sequence length of the PLM to retain as much contextual information as possible when they are fed into the PLM, we sequentially truncate the last tokens of the maximum-length utterances of the conversation during preprocessing until the maximum sequence length requirement of the PLM is met. The input sequence X^t is then fed into the

PLM to obtain a sequential representation of the entire conversation:

$$\mathbf{I}^t = PLM(X^t), \quad (11)$$

where $\mathbf{I}^t \in \mathbb{R}^{L \times d_t}$, L is the length of the input sequence and d_t is the hidden dimension of the PLM. To obtain the sequence representation of each utterance, we make a weighted average of the sequence representations of the tokens of each conversation in \mathbf{I}^t to obtain the sequence representation of each utterance $\mathbf{E}^t \in \mathbb{R}^{N \times d_t}$, where N denotes the number of utterances of that conversation. We selected Speaker-Aware RoBERTa (SA-RoBERTa) (Gu et al., 2020) as the PLM.

Audio : After resampling the audio to 16khz, we input it into an audio feature extraction model (AFE) to get a sequential representation of the audio modality of the conversation:

$$\mathbf{E}^a = AFE(X^a), \quad (12)$$

where $\mathbf{E}^a \in \mathbb{R}^{n \times d_a}$, and d_a is the hidden layer dimension of the audio feature extraction model. We choose Wav2Vec2.0 (Baevski et al., 2020) as the audio feature extraction model.

Video: We first sample the video at equal intervals as a sequence of images over several frames to obtain the input sequence \mathbf{X}^v , $\mathbf{X}^v \in \mathbb{R}^{F \times d_f \times d_f}$ of the video modality, where F is the number of sampled frames and d_f is the size of the picture. The image sequences are then fed into the video feature extraction model (VFE) to get a sequence representation of the video modalities:

$$\mathbf{E}^v = VFE(\mathbf{X}^v), \quad (13)$$

where $\mathbf{E}^v \in \mathbb{R}^{n \times d_v}$ and d_v is the hidden layer dimension of the video feature extraction model. We select the pre-trained DenseNet (Huang et al., 2017) as the audio feature extraction model.

C.2 Pseudo-code of Graph Construction

The pseudo-code of the graph construction process is shown in Algorithm C.2.

C.3 An Example of the Graph construction

If $K = 1$, the graph constructed for the conversation in Figure 1 is shown in Figure 6.

Table 4: Performance comparison of different methods for conversations with varying numbers of utterances. The best results and the second best results are in bold and underlined, respectively.

Method	ECF				MECAD			
	num_utt ≤ 10		num_utt > 10		num_utt ≤ 10		num_utt > 10	
	6 Avg.	4 Avg.	6 Avg.	4 Avg.	6 Avg.	4 Avg.	6 Avg.	4 Avg.
RankCP	31.50	33.29	29.34	31.88	27.19	29.23	25.11	27.13
SHARK	33.68	35.57	31.49	33.17	28.32	30.75	27.01	29.41
GPT-4o	30.08	31.56	28.42	29.36	26.34	27.82	27.79	28.88
M ³ HG (T)	<u>39.18</u>	<u>41.25</u>	<u>36.18</u>	<u>38.98</u>	<u>31.95</u>	<u>33.40</u>	<u>29.94</u>	<u>31.90</u>
M ³ HG (T, A, V)	41.95	40.42	38.67	41.09	33.76	35.21	32.10	34.12

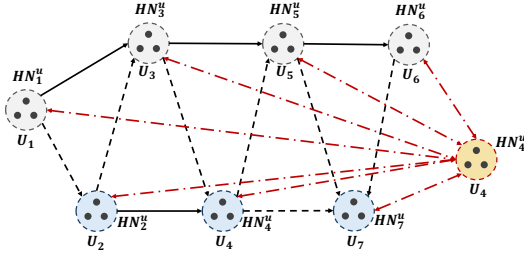


Figure 6: Super-Node-based edges and relations constructed from a conversation in MECAD with $K = 1$. The utterance Super-Nodes of the two speakers are shown in gray and blue, respectively. The black solid and dashed lines denote the Super-Edges between the same speaker and different speakers, respectively, and the red dotted lines denote the Super-Edges between the utterance Super-Nodes and the conversation Super-Nodes.

D Implement Details of the Experiment

For the ECF dataset, we use the pre-trained RoBERTa-large² model to initialize the feature extraction parameters of the text modality. For audio modality, we use the wav2vec2-base-960h³ model and for video modality we use the DenseNet (Huang et al., 2017) model. For the MECAD dataset, we use the chinese-roberta-wwm-ext-large⁴ model for the initialization of textual modal features, the wav2vec2-large-chinese-zh-cn⁵ model for the extraction of audio modal features, and the DenseNet model is also applied to the video modal. In our experiments, none of the parameters of the PLM were frozen. During the construction of the graph, we set the hyperparameter K to 3.

²<https://huggingface.co/FacebookAI/roberta-large>

³<https://huggingface.co/facebook/wav2vec2-base-960h>

⁴<https://huggingface.co/hfl/chinese-roberta-wwm-ext-large>

⁵<https://huggingface.co/wbbbbb/wav2vec2-large-chinese-zh-cn>

During training, we use the AdamW (Loshchilov, 2017) optimizer with batch size and learning rate set to 16 and 5e-6, respectively, and perform a parameter update after every two mini-batches. Our model is trained for 50 epochs on the training set, and the checkpoints corresponding to the highest values of the weighted average F1 scores of the six emotions on the validation set are used as the results of the test set.

E Supplementary Experimental Results of M³HG

E.1 Ablation Study

Effect of different modules. We conduct ablation studies to verify the effectiveness of different modules in M³HG on the two datasets using 6 Avg and 4 Avg scores. As shown in Table 6, *w/o N^e & N^c* indicates no use of emotional and causal context nodes in graph construction. Consequently, emotion-cause pair prediction is performed directly based on the features of each utterance node. *w/o inter-fusion* and *w/o intra-fusion* denote the absence of inter-utterance and intra-utterance multi-modal fusion, respectively, during multi-scale semantic information fusion. Our model outperforms the state-of-the-art baselines even without utilizing the previous three mechanisms. Specifically, the performance of M³HG degrades on both ECF and MECAD datasets when removing the emotional and causal context nodes, demonstrating the necessity of explicitly modeling the emotion and cause-related contexts. Moreover, removing both intra-utterance and inter-utterance semantic fusion results in a drop in the model’s performance, the former of which causes a more significant degradation. It highlights the importance of effectively fusing semantic information at different scales within heterogeneous graphs, particularly within individual utterances.

Table 5: Performance comparison of different methods on four subtasks. The best results and the second best results are in bold and underlined, respectively.

Dataset	Method	EP			ER			CE			EC		
		P	R	F1	P	R	F1	P	R	F1	P	R	F1
ECF	SHARK	<u>59.00</u>	61.21	<u>60.74</u>	<u>40.34</u>	45.65	<u>42.83</u>	<u>69.25</u>	66.13	67.64	<u>50.12</u>	46.31	<u>48.14</u>
	GPT-4o(5-shots)	45.17	78.62	57.37	36.42	42.70	36.76	57.02	84.85	<u>68.21</u>	32.90	61.13	42.78
	M ³ HG (T)	71.36	<u>75.11</u>	73.19	52.24	<u>45.63</u>	46.60	72.32	<u>68.40</u>	70.30	58.03	<u>52.05</u>	54.88
MECAD	SHARK	69.30	67.02	<u>68.14</u>	39.38	<u>36.93</u>	<u>38.12</u>	64.18	66.36	65.24	<u>49.02</u>	<u>42.87</u>	<u>45.74</u>
	GPT-4o(5-shots)	<u>71.41</u>	63.69	67.33	38.69	36.59	34.22	<u>65.03</u>	<u>69.26</u>	<u>67.08</u>	39.68	41.77	40.70
	M ³ HG (T)	72.34	67.84	70.02	43.35	40.98	41.66	66.12	70.24	68.12	54.82	46.42	50.27

Table 6: Ablation results.

Dataset	Model	6 Avg.	4 Avg.
ECF	M ³ HG	40.07	41.96
	<i>w/o</i> all modules	36.81($\downarrow 3.26$)	38.57($\downarrow 3.39$)
	<i>w/o</i> N^e & N^c	38.13($\downarrow 1.94$)	40.11($\downarrow 1.85$)
	<i>w/o</i> inter-fusion	39.56($\downarrow 0.51$)	41.14($\downarrow 0.82$)
	<i>w/o</i> intra-fusion	39.12($\downarrow 0.95$)	40.86($\downarrow 1.10$)
MECAD	M ³ HG	32.82	34.59
	<i>w/o</i> all modules	30.37($\downarrow 2.45$)	32.27($\downarrow 2.32$)
	<i>w/o</i> N^e & N^c	30.94($\downarrow 1.88$)	32.79($\downarrow 1.80$)
	<i>w/o</i> inter-fusion	32.57($\downarrow 0.25$)	33.91($\downarrow 0.68$)
	<i>w/o</i> intra-fusion	32.16($\downarrow 0.66$)	33.33($\downarrow 1.26$)

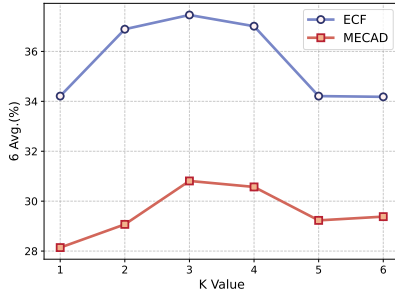


Figure 7: Results of M³HG with various K values.

Effect of the hyperparameter K . The hyperparameter K is closely related to the spatio-temporal complexity of the M³HG’s graph construction. We vary the size of K (ranging from 1 to 6) to test its effect, and the result of the M³HG on both datasets is shown in Figure 7. The performance of M³HG on both datasets initially improves with increasing K and then declines, with the best performance observed at $K = 3$.

E.2 In-Depth Analysis

The impact of conversation length. To evaluate the performance of M³HG in handling longer conversations, we present a comparison of the performance of M³HG and other baseline models across conversations of varying lengths, as shown in Table 4. We observe that M³HG outperforms all baseline models in scenarios involving conversations with more than 10 utterances, which account for


Table 7: Performance comparison of different methods for conversations in which cause utterance appears after emotion utterances. The best results and the second best results are in bold and underlined, respectively.

Method	ECF		MECAD	
	6 Avg.	4 Avg.	6 Avg.	4 Avg.
SHARK	29.15	30.54	25.49	27.51
GPT-4o	28.21	29.45	26.42	27.76
M ³ HG (T)	<u>35.48</u>	<u>37.01</u>	<u>29.18</u>	<u>30.93</u>
M ³ HG (T, A, V)	38.25	40.01	31.27	33.09

42.65% and 43.38% of all conversations in the ECF and MECAD datasets, respectively. In long conversations, baseline models, including GPT-4o, fail to effectively extract global contextual information, thereby missing a number of triplets. Our model, through semantic fusion at different scales within multimodal heterogeneous graphs, effectively captures more triplets by extracting contextual information from long conversations.

Model performance when cause utterance appears after emotion utterances. A key challenge in the MECATEC task is when the cause of a speaker’s emotion is revealed later in the conversation, requiring the model to effectively capture and interpret the global context of the conversation. To further emphasize the superior performance of M³HG in handling cases where the cause utterance appears after the emotion utterance, we identified and filtered all such conversations from the ECF and MECAD datasets. The performance of M³HG, compared with two other representative models, is shown in Table 7. M³HG demonstrates superior performance, while SHARK suffers a greater performance drop compared to M³HG. Although the performance drop for GPT-4o (5-shots) is less pronounced, its overall performance remains unsatisfactory.

Model performance on four subtasks. To evaluate M³HG’s performance more comprehensively,

U_1		Yunru Chen: 干嘛这样看着我啊? (Why are you looking at me like that?)
U_2		Junjie Mo: 没事啦, 不用客气。 (It's nothing, no need to thank me.)
U_3		Yunru Chen: 你是不是跟李子维一样, 觉得我说的那些话, 都是乱编的? (Are you like Ziwei Li, thinking that what I said was all made up?)
U_4		Junjie Mo: 我相信你说的都是真的啊, 在你的梦里, 真的有那么一个人, 你很喜欢他, 他也很喜欢你, 而且... (I believe what you said is true. In your dream, there was really someone you liked a lot, and he liked you too, and...)
U_5		Junjie Mo: 没事啦 (It's nothing.)
U_6		Yunru Chen: 而且什么, 你说啊? (And what? Tell me!)
U_7		Junjie Mo: 也许比起李子维, 我更希望你喜欢的, 只是你梦里那个王詮胜。 (Maybe, compared to Li Ziwei, I wish you'd like only the Wang Quansheng in your dream.)
Ground Truth		(Surprise,1,1), (Sad,3,3), (Sad,5,7), (Surprise,6,4), (sad,7,7)
SHARK		(Surprise,1,1), (Sad,3,3), (Surprise,6,6), (Sad,7,6)
GPT-4o (5-shots)		(Anger,3,1), (Surprise,6,4), (Sad,7,7)
M ³ HG		(Surprise,1,1), (Sad,3,3), (Sad,5,7), (Anger,6,5)

U_1		Zongming Tan: 怎么了? (What's going on?)
U_2		Di An: 我也不知道, 总觉得有人在跟着我。 (I don't know. I always feel like someone's following me.)
U_3		Zongming Tan: 你刚回来不久, 上海本身就没几个朋友, 谁会跟着你。 (You just came back not long ago, and you don't have many friends in Shanghai itself, who would follow you.)
U_4		Di An: 我也觉得奇怪, 加上今天已经好几次了, 总觉得有人在跟着我, 一回头, 又什么都没有, 你说, 会不会是我自己的幻觉, 还是? (I also think it's strange, plus it's been several times today, I always feel that someone is following me, and when I turn around, there's nothing.)
U_5		Zongming Tan: 安迪, 别胡思乱想, 可能就是工作太辛苦, 太累了。 (Andy, don't get any ideas, it's probably just a case of working too hard and being too tired.)
U_6		Di An: 可能吧, 也许是我今天没有吃早餐, 低血糖了, 所以才有幻觉。 (Maybe, maybe I'm hallucinating because I didn't eat breakfast today and I'm low on blood sugar.)
Ground Truth		(Fear,2,2), (Fear,4,2), (Fear,4,4)
SHARK		(Fear,4,4)
GPT-4o (5-shots)		(Surprise,1,1), (Fear,2,2), (Sad,4,4), (Anger,5,5), (Sad,6,6)
M ³ HG (T)		(Fear,2,2), (Fear,4,3), (Fear,4,4)
M ³ HG (T+A+V)		(Fear,4,2), (Fear,4,4), (Sad,6,6)

Figure 8: Comparison of utter-cause-emotion triplet on two test samples.

we define the following four subtasks:

- **Emotion Extraction (EP):** Predict whether an utterance expresses an emotion (binary classification), same as SHARK.
- **Emotion Recognition (ER):** Predict the emotion category of an utterance (multi-class classification).
- **Cause Extraction (CE):** Predict whether an utterance is a cause utterance (binary classification), same as SHARK.
- **Emotion-Cause Pair Extraction (EC):** Predict whether two utterances of a conversation form an emotion-cause pair (binary classification).

Table 5 demonstrates the performance comparison between M³HG and other SOTA models across the four subtasks. For the EP subtask, M³HG performs the best across both datasets. It is worth noting that GPT-4o (5-shots) achieves a high recall on the ECF dataset. This phenomenon can be attributed to the more pronounced label sparsity in the ECF dataset compared to MECAD. As a result, GPT-4o (5-shots) frequently predicts that an utterance carries emotion, leading to a higher recall. For the ER subtask, M³HG achieves the best results across both datasets. This demonstrates M³HG’s ability to effectively extract the emotional context embedded in utterances. For the CE subtask, M³HG performs best, demonstrating the importance of integrating the cause prediction subtask

into the model during training. For the EC subtask, GPT-4o (5-shots) similarly exhibits high recall on the ECF dataset. This is due to the severe label sparsity problem in the ECF dataset, compared to MECAD, which leads GPT-4o to predict as many emotion-cause pairs as possible.

E.3 Case Study

To demonstrate the superiority and limitations of M³HG, we present a case study that compares the prediction results of M³HG with those of two other representative models (i.e. SHARK, GPT-4o (5-shots)), using two sample conversations from the MECAD dataset. As shown in Figure 8, the first test sample demonstrates that M³HG outperforms the other models in prediction accuracy, while GPT-4o exhibits the poorest performance. This can be attributed to M³HG’s use of a multimodal heterogeneous graph and a specially designed conversation super-node, which effectively captures global contextual information. These features enable M³HG to more accurately handle scenarios where the cause utterance appears after the emotion utterance.

In the second sample, M³HG (T+A+V) is less effective than M³HG (T) in predicting Utterance 6 as “Sad” and Utterance 2 as “Neutral”. This is because the combination of text and context in Utterance 2 conveys the speaker’s worried and fearful mood, while the video and audio signals suggest a

Table 8: An example of prompt for ChatGPT.

Instruction	<p>You are an expert in sentiment analysis and identification of emotional causes. I will give you a conversation between multiple speakers. You are required to extract the utter-cause-emotion triplet for a given utterance. First, infer the emotion label for the utterance (select one from: Anger, Disgust, Fear, Joy, Sadness, Surprise or Neutral). Then, identify the index(es) of the cause utterance(s) that triggered this emotion (the index should represent the utterance(s) from the conversation that caused the emotion, and it must be non-negative. Multiple indices should be separated by commas). If the predicted emotion is Neutral, there is no corresponding cause utterance. The output should follow the format: emotion label, cause utterance indices. Examples of the expected output format: Example 1: happy,3. Example 2: sad,3,4,5. Example 3: neutral.</p>
Input	<p>Input Conversation :</p> <p>{ 1. Fang Sijin: First, change your clothes, then head to this address. A decoration company will be coming over shortly. You'll need to supervise their work and see how you can help. }</p> <p>{ 2. Zhu Shanshan: Wait, am I really responsible for this? I don't know anything about decoration. }</p> <p>{ 3. Fang Sijin: You've been handing out flyers for two days now. Have you gotten any interested customers? }</p> <p>{ 4. Zhu Shanshan: But you only told me to distribute the flyers; you never ask for phone numbers! }</p> <p>Candidate Utterances:</p> <p>{ 1. Fang Sijin: First, change your clothes, then head to this address. A decoration company will be coming over shortly. You'll need to supervise their work and see how you can help. }</p> <p>{ 2. Zhu Shanshan: Wait, am I really responsible for this? I don't know anything about decoration. }</p> <p>Target Utterance:</p> <p>{ 2. Zhu Shanshan: Wait, am I really responsible for this? I don't know anything about decoration. }</p> <p>Target emotion labels and cause index(es):</p> <p>[Surprise, 1]</p> <p>Input Conversation :</p> <p>.....</p> <p>Candidate Utterances:</p> <p>.....</p> <p>Target Utterance:</p> <p>.....</p> <p>Target emotion labels and cause index(es):</p> <p>.....</p>
Output	<p>output example [Happy, 1, 2]</p>

calmer demeanor. This discrepancy likely caused M³HG (T+A+V) to mispredict the emotions in this case. Nevertheless, M³HG still outperforms all other baseline models, demonstrating its robustness and superior predictive capability even under challenging conditions.

F Prompt Design for ChatGPT

We use the GPT-4o model of OpenAI public API (version up to May 13, 2024) and design a prompt elaborately to test the performance on the MECTEC task. The prompt (i.e., the input of ChatGPT) includes three parts:

- **Instruction.** We use instructions to guide the ChatGPT on what it needs to do. Our instruction is as follows:

You are an expert in sentiment analysis and identification of emotional causes. I will give you a conversation between two or more speakers. You need to extract the utter-cause-emotion triplet of the given utterance.

Meanwhile, we provide a detailed description of the output formats required for ChatGPT, as illustrated in Table 8.

- **Demonstrations** We achieve the few-shot in-context learning of ChatGPT by adding demonstrations. We use the 5-shot in-context

learning due to the limitations of the input length. Each demonstration includes a conversation as input and a target utterance as the target for prediction.

Except for the aforementioned two parts, we also need to describe the conversations to be predicted and the corresponding target utterance. An example is shown in Table 8.