

Emotion–Cause Pair Extraction in Conversations via Semantic Decoupling and Graph Alignment

Anonymous ACL submission

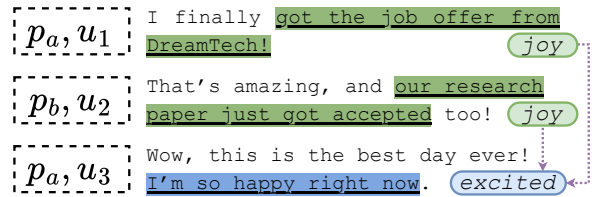
Abstract

Emotion-Cause Pair Extraction in Conversations (ECPEC) aims to identify the set of causal relations between emotion utterances and their triggering causes within a dialogue. Most existing approaches formulate ECPEC as an independent pairwise classification task, overlooking the distinct semantics of emotion diffusion and cause explanation, and failing to capture globally consistent many-to-many conversational causality. To address these limitations, we revisit ECPEC from a semantic perspective and seek to disentangle emotion-oriented semantics from cause-oriented semantics, mapping them into two complementary representation spaces to better capture their distinct conversational roles. Building on this semantic decoupling, we naturally formulate ECPEC as a global alignment problem between the emotion-side and cause-side representations, and employ optimal transport to enable many-to-many and globally consistent emotion-cause matching. Based on this perspective, we propose a unified framework SCALE that instantiates the above semantic decoupling and alignment principle within a shared conversational structure. Extensive experiments on several benchmark datasets demonstrate that SCALE consistently achieves state-of-the-art performance.¹

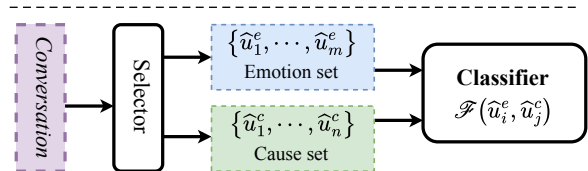
1 Introduction

Emotion-Cause Pair Extraction in Conversations (ECPEC) aims to identify the set of causal relations between emotion utterances and their triggering causes within a dialogue, which often exhibit complex and many-to-many dependencies. Figure 1a illustrates a representative example, where the emotion utterance u_3 is jointly triggered by multiple preceding utterances u_1 and u_2 . Unlike emotion recognition in conversation (ERC) (Wang et al., 2024b; Fu et al., 2023; Majumder et al., 2019),

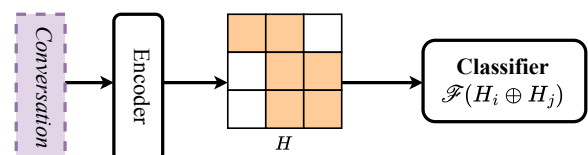
¹The code and data are available at <https://anonymous.4open.science/r/r9f3k2/>.



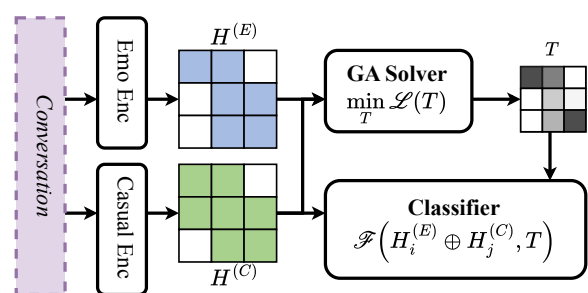
(a) An example of ECPEC task.



(b) Select-then-pair paradigm.



(c) Embed-then-pair paradigm.



(d) Our proposal.

Figure 1: Comparison between the existing ECPEC paradigm (b-c) and SCALE (d).

which focuses on assigning discrete emotion labels to individual utterances (Fu et al., 2021; Hu et al., 2021), ECPEC provides a causal perspective for dialogue understanding, enabling more fine-grained analysis of emotional dynamics. As highlighted by Poria et al. (2021), ECPEC has broad applicability across multiple domains, including dialogue systems (Rashkin et al., 2019; Zhong et al., 2020),

041
042
043
044
045
046
047
048

049 conversational recommendation (Liang et al., 2024;
050 Li et al., 2018), mental health analysis (Cambria
051 et al., 2018; Pontiki et al., 2016), and social me-
052 dia opinion mining (Alexander Pak and Patrick
053 Paroubek, 2010; Liu, 2022).

054 Early studies on ECPEC predominantly followed
055 the *select-then-pair* paradigm (Ding et al., 2020;
056 Wang et al., 2023a), which independently iden-
057 tifies candidate emotion utterances and cause ut-
058 terances before pairing them through heuristic or
059 classifier-based matching, as shown in 1c. While
060 intuitive and easy to integrate with existing ERC
061 models (Gao et al., 2023; Nguyen et al., 2024;
062 Ghosal et al., 2019), this pipeline is prone to error
063 propagation during candidate selection and fails to
064 fully exploit contextualized utterance representa-
065 tions. To alleviate these issues, subsequent studies
066 shifted towards the *embed-then-pair* paradigm (An
067 et al., 2023; Li et al., 2023; Jeong and Bak, 2023;
068 Wang et al., 2024a), where utterance embeddings
069 are directly concatenated and classified as emotion-
070 cause pairs in an end-to-end manner. Although this
071 paradigm better leverages utterance-level seman-
072 tics and avoids explicit candidate construction, it
073 still treats emotion-cause inference as a collection
074 of independent pairwise decisions.

075 Despite their procedural differences, existing
076 ECPEC approaches share two fundamental *limi-*
077 *tations*. **L1)** Most methods encode emotion-related
078 and cause-related information within a unified rep-
079 resentation space or interaction structure, implic-
080 itly assuming that emotion diffusion and cause ex-
081 planation follow homogeneous relational patterns.
082 However, in real conversations, emotional states
083 tend to propagate through contextual and speaker-
084 dependent dynamics, whereas causes are grounded
085 in explanatory and often asymmetric dependencies.
086 Conflating these distinct semantics obscures their
087 respective roles in conversational causality. **L2)**
088 Existing methods typically formulate ECPEC as
089 independent one-to-one pair classification with bi-
090 nary judgments. Such pairwise formulations are
091 inherently inadequate for modeling globally con-
092 sistent many-to-many causal structures, where mul-
093 tiple interdependent causes may jointly trigger an
094 emotion and a single cause may influence multiple
095 emotional outcomes.

096 To address these limitations, we revisit ECPEC
097 from a semantic perspective and argue that emotion
098 diffusion and cause explanation, while grounded in
099 the same conversational structure, should be charac-
100 terized by different semantic focuses. Rather than

101 duplicating dialogue structures or enforcing task-
102 level separation, we seek to disentangle emotion-
103 oriented and cause-oriented semantics by map-
104 ping them into two complementary representation
105 spaces induced from a shared conversation graph.
106 Building on this semantic decoupling, we naturally
107 formulate ECPEC as a global alignment problem
108 between emotion-side and cause-side representa-
109 tions, which enables holistic reasoning over many-
110 to-many emotion-cause relations. Based on this
111 perspective, we propose **SCALE** (**S**emantic **C**ausal
112 **A**lignment for **E**CPEC), a unified framework that
113 instantiates semantic decoupling and global align-
114 ment within conversational contexts. Extensive ex-
115 periments on multiple benchmark datasets demon-
116 strate that SCALE consistently outperforms exist-
117 ing state-of-the-art approaches. Overall, the main
118 contributions of this work are summarized as fol-
119 lows:

- We revisit ECPEC from a semantic perspec- 120
121 tive and highlight the necessity of disentan-
122 gling emotion diffusion and cause explanation
123 while preserving shared conversational struc-
124 ture.
- We propose SCALE, a unified framework that 125
126 induces emotion-side and cause-side repre-
127 sentations from a shared conversation graph
128 and formulates ECPEC as a global alignment
129 problem to support many-to-many and glob-
130 ally consistent inference.
- Extensive experiments on several public 131
132 ECPEC benchmarks demonstrate that SCALE
133 consistently achieves state-of-the-art perfor-
134 mance.

135 2 Methodology

136 Formally, given a conversation $\mathcal{C}_i =$ 136
137 $\{(u_1, p_{\pi(u_1)}), \dots, (u_N, p_{\pi(u_N)})\}$ consisting
138 of N utterances, each utterance u_j is associ-
139 ated with a speaker $p_{\pi(u_j)}$, where π denotes
140 a mapping from an utterance u_i to the index
141 of its corresponding speaker. In ECPEC, the
142 goal is to identify a set of emotion-cause pairs
143 $\mathcal{P} = \{(u_e, u_c) \mid u_e \text{ is caused by } u_c\}$ that char-
144 acterizes the underlying conversational causality.
145 To address the limitations of existing approaches,
146 we propose a framework termed SCALE that
147 provides a unified solution for the ECPEC task.
148 An overview of the proposed SCALE is illustrated
149 in Figure 2. Specifically, each conversation is first

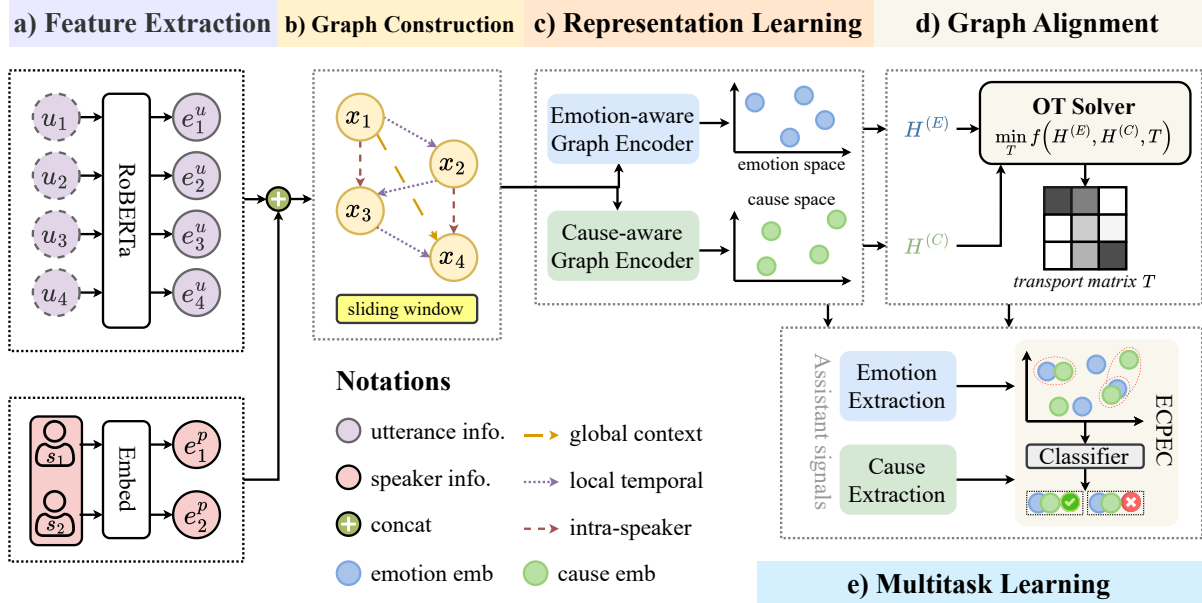


Figure 2: Overall architecture of our SCALE.

150 encoded into utterance-level representations and
 151 organized as a conversation graph (§2.1). Then,
 152 SCALE induces two complementary semantic
 153 views of the dialogue, namely emotion-oriented
 154 and cause-oriented representations, by applying
 155 semantic-specific graph encoding mechanisms
 156 (§2.2). To explicitly model the correspondence
 157 between emotional utterances and their underlying
 158 causes, SCALE further formulates emotion-cause
 159 inference as a global alignment problem between
 160 the emotion-side and cause-side representations,
 161 which is solved via an optimal transport framework
 162 to enable many-to-many and globally consistent
 163 matching (§2.3). All components are optimized
 164 under a unified learning objective, where emotion
 165 extraction and cause extraction are introduced
 166 as auxiliary supervision signals to facilitate
 167 representation learning and ultimately improve
 168 ECPEC performance (§2.4).

169 2.1 Encode Conversation as a Graph

170 To establish relationships between utterances while
 171 capturing both inter- and intra-speaker depend-
 172 encies (Ghosal et al., 2019; Li et al., 2023; Gao
 173 et al., 2023), we represent each conversation as
 174 a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{A})$. Each node $v_i \in \mathcal{V}$
 175 corresponds to an utterance u_i , edges $e_{ij} \in \mathcal{E}$
 176 are constructed according to three types of dependencies,
 177 and $\mathbf{A} \in \mathbb{R}^{N \times N}$ is the adjacency matrix.

178 **Nodes.** Each utterance u_i is represented as a
 179 node v_i , initialized with a contextual embedding

180 $\mathbf{x}_i^u \in \mathbb{R}^{d_u}$ obtained from a pretrained RoBERTa
 181 model (Liu et al., 2019). To incorporate speaker in-
 182 formation, the corresponding speaker $p_{\pi(u_i)}$ is first
 183 represented as a one-hot vector and then projected
 184 into a speaker embedding $\mathbf{x}_{\pi(u_i)}^s \in \mathbb{R}^{d_s}$ through
 185 a learnable embedding layer. The final speaker-
 186 aware utterance-level feature is defined as:

$$187 \mathbf{x}_i = \mathbf{x}_i^u \oplus \mathbf{x}_{\pi(u_i)}^s, \quad (1)$$

188 where $\mathbf{x}_i \in \mathbb{R}^{d_h}$ is the representation of node v_i , \oplus
 189 denotes the concatenation operation.

190 **Edges.** We establish three types of relationships
 191 between utterance nodes to capture both global and
 192 local contextual dependencies in the conversation
 193 graph \mathcal{G} . The global contextual edge (v_i, v_j)
 194 captures long-range semantic dependencies between
 195 utterance nodes v_i and v_j . Such an edge exists
 196 iff $\cos(\mathbf{x}_i, \mathbf{x}_j) + 1 > \tau_s$, where \mathbf{x}_i and \mathbf{x}_j
 197 denote the node features of v_i and v_j , respectively,
 198 and τ_s is a similarity threshold. The local contextual
 199 edge models short-range emotional dynamics. An
 200 edge (v_i, v_j) exists if $|i - j| \leq W$, where W
 201 denotes the size of the sliding temporal window. The
 202 intra-speaker edge captures speaker-specific emo-
 203 tional consistency. We connect nodes v_i and v_j
 204 iff $\pi(u_i) = \pi(u_j)$ and $i \neq j$.

205 **Edge weights.** We initialize edge weights accord-
 206 ing to the corresponding dependency type. For
 207 global contextual edges, we initialize weights based

on utterance-level semantic similarity:

$$e_{ij} = \frac{\cos(\mathbf{x}_i, \mathbf{x}_j) + 1}{2}. \quad (2)$$

For local temporal edges, weights decay exponentially with conversational distance:

$$e_{ij} = \exp(-|i - j|/\tau_e), \quad (3)$$

where τ_e controls temporal sensitivity. For intra-speaker edges, weights reflect speaker-specific emotional consistency across turns:

$$e_{ij} = \frac{\exp(-|i - j|/\tau_e) + 1}{2}. \quad (4)$$

All initialized edge weights are integrated into a single weighted adjacency matrix \mathbf{A} , which is treated as learnable and jointly optimized with other model parameters.

2.2 Graph Representation Learning

Given the conversation graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{A})$ constructed for emotion-side and cause-side modeling, respectively, along with the same feature matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times d_h}$, we employ two independent graph encoders, GNN_E and GNN_C , to learn two sets of node representations in different semantic spaces:

$$\mathbf{H}^{(S)} = \text{GNN}_S(\mathcal{G}, \mathbf{X}), \quad (5)$$

where $S \in \{E, C\}$ denotes the emotion or cause semantic space, $\mathbf{H}^{(E)} \in \mathbb{R}^{N \times d_h}$ and $\mathbf{H}^{(C)} \in \mathbb{R}^{N \times d_h}$ denote the resulting emotion-aware and cause-aware node representations. At each layer in encoder, node representations $\mathbf{h}^{(S)} \in \mathbf{H}^{(S)}$ are updated via attention-weighted message passing:

$$\begin{aligned} \psi_{ij}^{(S)} &= \phi^{(S)}(\mathbf{h}_i^{(S)}, \mathbf{h}_j^{(S)}) A_{ij}, \\ \alpha_{ij}^{(S)} &= \text{softmax}_{j \in \mathcal{N}(i)}(\psi_{ij}^{(S)}), \\ \mathbf{h}_i^{(S)} &= \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^{(S)} \mathbf{W}^{(S)} \mathbf{h}_j^{(S)}, \end{aligned} \quad (6)$$

where $\mathcal{N}(i)$ denotes the neighborhood of node v_i in \mathcal{G} , $\alpha_{ij}^{(S)}$ is the normalized attention coefficient, $\mathbf{W}^{(S)} \in \mathbb{R}^{d_h \times d_h}$ is a learnable matrix, and $\phi^{(S)}(\cdot)$ is a learnable scoring function. Importantly, the attention coefficient $\alpha_{ij}^{(S)}$ learned by each encoder can be viewed as a semantic-specific edge weight between v_i and v_j over the shared graph. Specifically, we define a task-specific weighted adjacency matrix $\mathbf{A}^{(S)} \in \mathbb{R}^{N \times N}$ as

$$A_{ij}^{(S)} = \alpha_{ij}^{(S)}, \quad (7)$$

which captures the relative importance of edge (v_i, v_j) under semantic space S . As a result, the emotion-aware and cause-aware encoders implicitly induce two refined adjacency structures, denoted as $\mathbf{A}^{(E)}$ and $\mathbf{A}^{(C)}$, respectively.

2.3 Graph Alignment via Optimal Transport

To capture the many-to-many relations inherent in emotion-cause pair extraction, we formulate ECPEC as a global alignment problem between emotion-aware and cause-aware representations. Formally, given the emotion representations $\mathbf{H}^{(E)} = \{\mathbf{h}_i^{(E)}\}_{i=1}^N$ and the cause representations $\mathbf{H}^{(C)} = \{\mathbf{h}_i^{(C)}\}_{i=1}^N$, our goal is to learn a transport plan $\mathbf{T} \in \mathbb{R}^{N \times N}$, where each entry $T_{ij} \geq 0$ indicates the soft correspondence strength between the i -th emotion representation $\mathbf{h}_i^{(E)}$ and the j -th cause representation $\mathbf{h}_j^{(C)}$. Unlike independent pairwise scoring, the alignment is learned globally over the entire matrix \mathbf{T} , such that multiple causes can be jointly associated with the same emotion and the correspondences across different emotion-cause pairs are mutually constrained.

Alignment objective. We seek to learn the alignment matrix \mathbf{T} by minimizing the following objective with respect to \mathbf{T} :

$$\begin{aligned} \min_{\mathbf{T} \geq 0} \mathcal{L}(\mathbf{T}) &= \alpha \langle \mathbf{C}_{\text{attr}}, \mathbf{T} \rangle \\ &+ (1 - \alpha) \mathcal{L}_{\text{struct}}(\mathbf{T}) \end{aligned} \quad (8)$$

where $\mathbf{C}_{\text{attr}} \in \mathbb{R}^{N \times N}$ denotes the attribute-level cost matrix, $\mathcal{L}_{\text{struct}}(\mathbf{T})$ denotes a structure consistency term. \mathbf{C}_{attr} measures semantic compatibility between emotion representation $\mathbf{h}_i^{(E)} \in \mathbf{H}^{(E)}$ and cause representation $\mathbf{h}_j^{(C)} \in \mathbf{H}^{(C)}$, with each entry defined as:

$$\mathbf{C}_{\text{attr}}(i, j) = 1 - \cos(\mathbf{h}_i^{(E)}, \mathbf{h}_j^{(C)}), \quad (9)$$

where smaller values indicate higher semantic affinity. The structure-level term $\mathcal{L}_{\text{struct}}(\mathbf{T})$ is defined as a structure consistency loss that measures the discrepancy between relational patterns encoded in the emotion-side and cause-side dialogue graphs. Specifically, it is formulated as:

$$\mathcal{L}_{\text{struct}}(\mathbf{T}) = \sum_{i,k,j,l} |A_{ik}^{(E)} - A_{jl}^{(C)}|^2 T_{ij} T_{kl}. \quad (10)$$

Optimization. The resulting objective $\mathcal{L}(\mathbf{T})$ is non-linear due to the quadratic structure consistency term $\mathcal{L}_{\text{struct}}(\mathbf{T})$. To efficiently minimize it in a differentiable manner, we adopt an entropy-regularized Sinkhorn scheme based on the standard fused Gromov-Wasserstein optimization strategy. Starting from a uniform initialization $\mathbf{T}^{(0)}$, we iteratively linearize the structure consistency term around the current solution and solve a sequence of entropic optimal transport subproblems. At iteration t , the linearized approximation of $\mathcal{L}_{\text{struct}}(\mathbf{T})$ induces an effective structure-aware cost matrix:

$$\mathbf{C}_{\text{struct}}^{(t)}(i, j) = \sum_{k, l} \left| \mathbf{A}_{ik}^{(E)} - \mathbf{A}_{jl}^{(C)} \right|^2 T_{kl}^{(t)}, \quad (11)$$

which leads to the following cost matrix used to construct a linear surrogate of $\mathcal{L}(\mathbf{T})$:

$$\mathbf{C}^{(t)} = \alpha \mathbf{C}_{\text{attr}} + (1 - \alpha) \mathbf{C}_{\text{struct}}^{(t)}. \quad (12)$$

The updated alignment matrix \mathbf{T} is iteratively updated via Sinkhorn normalization:

$$\mathbf{T}^{(t+1)} = \mathcal{S} \left(\exp(-\mathbf{C}^{(t)}/\varepsilon) \right), \quad (13)$$

where \mathcal{S} denotes the Sinkhorn operator with standard marginal constraints, and ε is the entropy regularization coefficient controlling the smoothness of the transport plan. After convergence, we apply a row-wise softmax with temperature τ_r to emphasize dominant alignments:

$$\tilde{\mathbf{T}} = \text{softmax}(\mathbf{T}/\tau_r). \quad (14)$$

The resulting $\tilde{\mathbf{T}}$ can be interpreted as a normalized alignment distribution, which reflects potential many-to-many associations between emotions and causes within each dialogue.

2.4 Multitask Joint Learning

SCALE adopts a multitask joint learning framework, where all task-specific objectives are optimized simultaneously with a shared encoder and task-specific prediction heads.

Emotion-Cause Pair Prediction. Given the row-normalized alignment matrix $\tilde{\mathbf{T}}$, each entry \tilde{T}_{ij} represents the global correspondence strength between an emotion utterance u_i and a candidate cause utterance u_j . To incorporate local discriminative evidence, we compute a pairwise matching score

by applying a lightweight classifier to the concatenated emotion-aware and cause-aware representations of each utterance pair:

$$s_{ij} = \mathcal{F}_{\text{ECPEC}} \left(\left[\mathbf{H}_i^{(E)}, \mathbf{H}_j^{(C)} \right] \right), \quad (15)$$

where $\mathcal{F}_{\text{ECPEC}}$ is implemented as a lightweight MLP with a sigmoid output layer. The final prediction score for emotion-cause pairs is obtained by combining global alignment and local evidence:

$$\hat{y}_{ij}^{(\text{ECPEC})} = \beta \tilde{T}_{ij} + (1 - \beta) s_{ij}, \quad (16)$$

where β controls the relative contribution of global correspondence and local evidence.

Emotion and Cause Extraction. We perform utterance-level emotion extraction (EE) and cause extraction (CE) using two parallel lightweight classifiers:

$$\begin{aligned} \hat{\mathbf{Y}}^{(E)} &= \mathcal{F}_{\text{EE}}(\mathbf{H}^{(E)}), \\ \hat{\mathbf{Y}}^{(C)} &= \mathcal{F}_{\text{CE}}(\mathbf{H}^{(C)}), \end{aligned} \quad (17)$$

where \mathcal{F}_{EE} and \mathcal{F}_{CE} are implemented as lightweight MLPs with softmax output layers.

Joint Optimization. The learning of SCALE is performed by minimizing \mathcal{L} :

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{\text{ECPEC}} + \lambda_{EE} \mathcal{L}_{\text{EE}} + \lambda_{CE} \mathcal{L}_{\text{CE}}, \\ \mathcal{L}_{\text{EE}} &= \text{CE}(\hat{\mathbf{Y}}^{(E)}, \mathbf{Y}^{(E)}), \\ \mathcal{L}_{\text{CE}} &= \text{CE}(\hat{\mathbf{Y}}^{(C)}, \mathbf{Y}^{(C)}), \end{aligned} \quad (18)$$

where \mathcal{L}_{EE} and \mathcal{L}_{CE} are cross entropy losses, λ_{EE} and λ_{CE} controlling their contributions. The ECPEC loss $\mathcal{L}_{\text{ECPEC}}$ is defined as:

$$\mathcal{L}_{\text{ECPEC}} = \mathcal{L}_{\text{pair}} + \lambda_{\text{OT}} \mathcal{L}_{\text{OT}}. \quad (19)$$

The pair-level supervision term is computed with binary cross-entropy over the predicted pair scores:

$$\mathcal{L}_{\text{pair}} = \text{BCE} \left(\hat{y}^{(\text{ECPEC})}, y \right), \quad (20)$$

where $y \in \{0, 1\}$ is the ground-truth label for the emotion-cause pair. The OT consistency regularizer encourages the local pair-wise prediction to match the OT-derived alignment score:

$$\mathcal{L}_{\text{OT}} = D_{\text{KL}} \left(\text{Bern}(s_{ij}) \parallel \text{Bern}(\tilde{T}_{ij}) \right). \quad (21)$$

Here D_{KL} denotes the Kullback–Leibler divergence, Bern is a Bernoulli distribution, s_{ij} is the local pair-wise prediction score in Equation (15), and $\tilde{T}_{ij} \in \mathbf{T}$ denotes the corresponding OT-derived alignment score.

Dataset	#Dlg.	#Utt.	#Pairs	Partition
RECCON-DD	1,106	11,104	5,861	75/5/20
RECCON-IE	16	665	1154	test only
ECF	1,374	13,619	9,794	70/10/20

Table 1: Dataset statistics.

3 Experiments

To comprehensively evaluate the proposed SCALE, we formulate the following *Research Questions* to guide our experiments:

RQ1: How does SCALE compare with existing state-of-the-art approaches on ECPEC datasets?

RQ2: How robust is SCALE in handling multi-cause scenarios compared to prior methods?

RQ3: How do key components of SCALE contribute to ECPEC performance?

RQ4: Can SCALE provide interpretable emotion-cause alignments and meaningful insights through qualitative analysis?

3.1 Experimental Setups

Datasets. We evaluate our method on three representative datasets: RECCON (Poria et al., 2021), which consists of two subsets, RECCON-DD and RECCON-IE, and ECF (Wang et al., 2023a). The dataset statistics are illustrated in Table 1, further details of these datasets are provided in Appendix B.

Baselines. We compare our method against several representative baselines, including RECCON (Poria et al., 2021) based on RoBERTa, MECPE-2steps (Wang et al., 2023a), PRG-MoE (Jeong and Bak, 2023), Joint-Xatt (Li et al., 2023), Joint-GCN (Li et al., 2023), MRC (Liu et al., 2023), and CENTER (Wang et al., 2024a). Detailed descriptions of these baselines are provided in Appendix C.

Metrics. Following previous work (Xia and Ding, 2019), we adopt F1-score (F1), Precision (P), and Recall (R) as evaluation metrics.

Implementation Details. We derive textual features for each utterance using a pretrained RoBERTa model. Unless otherwise specified, hyperparameters are set as follows: the temporal window size $W = 5$, utterance and speaker embedding dimensions $d_u = 768$ and $d_s = 50$, similarity and decay parameters $\tau_s = 0.5$, $\tau_e = 2.0$, and

$\tau_r = 1.0$, and weighting coefficients $\alpha = 0.8$, $\beta = 0.4$, $\varepsilon = 0.5$, $\lambda_{EE} = 0.2$, $\lambda_{CE} = 0.4$, and $\lambda_{OT} = 1.0$. Model training is conducted using the AdamW optimizer (Loshchilov and Hutter, 2017) with a learning rate of 1×10^{-4} . We employ a ReduceLROnPlateau scheduler and early stopping based on validation performance. All experiments are implemented in PyTorch (Paszke et al., 2019) and executed on a single NVIDIA RTX 4090 GPU. For a consistent comparison, all baselines are reimplemented based on their publicly released code or original descriptions, and trained under the same experimental settings.

3.2 Results and Discussion

3.2.1 Overall Performance (RQ1)

We evaluate all methods on the ECPEC task across three benchmarks, namely RECCON-DD, RECCON-IE, and ECF. As reported in Table 2, SCALE consistently achieves the highest recall and F1-score on all three datasets. On RECCON-DD, SCALE attains an F1 score of 58.83, outperforming the strongest baseline in terms of F1, PRG-MoE (57.26), by a relative improvement of +2.7%. On the smaller and cross-domain RECCON-IE dataset, SCALE achieves the best F1 score of 34.69, surpassing the strongest baseline (28.90) by a substantial relative gain of +20.0%, which highlights its strong generalization capability. Similarly, on ECF, SCALE delivers a notable relative improvement of +9.5% over MECPE-2steps (52.71) in terms of F1-score. We also observe that SCALE does not achieve the highest precision on most datasets. This behavior is consistent with the design of the soft optimal transport alignment, which encourages broader semantic matching between emotion and cause representations and therefore favors higher recall at the potential expense of precision.

3.2.2 Multi-Cause Study (RQ2)

As mentioned above, multi-cause scenarios introduce additional challenges for ECPEC. To evaluate model robustness under such settings, We therefore construct three multi-cause test subsets from RECCON-DD, RECCON-IE, and ECF, by selecting all dialogues in the original test splits where a single target emotion is annotated with two or more distinct causes. We then re-evaluate all models on these subsets to assess their robustness in capturing multiple causal triggers. As shown in Table 3, all models suffer from a notable performance drop, confirming the inherent difficulty of

	RECCON-DD			RECCON-IE			ECF		
	P	R	F1	P	R	F1	P	R	F1
RECCON	49.31	33.19	39.68	51.04	11.00	18.10	30.26	37.58	33.52
MECPE-2steps	49.34	47.37	48.34	27.31	6.30	10.24	57.64	48.72	<u>52.71</u>
PRG-MoE	58.95	<u>55.67</u>	<u>57.26</u>	<u>51.95</u>	20.02	<u>28.90</u>	47.11	55.27	50.86
Joint-Xatt	28.64	40.43	33.53	28.37	12.30	17.16	42.65	39.19	40.85
Joint-GCN	30.79	36.88	33.56	27.49	16.67	20.75	40.29	42.33	41.28
MRC	52.19	52.86	52.47	59.59	16.08	20.96	44.46	<u>57.65</u>	50.20
CENTER	47.39	46.88	47.13	34.92	<u>24.38</u>	28.71	47.90	43.65	44.75
SCALE	<u>56.31</u>	61.60	58.83	42.54	29.29	34.69	<u>55.01</u>	60.67	57.70

Table 2: Comparison results of ECPEC task.

	RECCON-DD	RECCON-IE	ECF
RECCON	28.27	8.76	23.75
MECPE-2steps	33.61	11.48	28.71
PRG-MoE	37.84	21.62	33.71
Joint-Xatt	23.18	4.62	25.15
Joint-GCN	23.73	5.08	26.96
MRC	33.39	3.22	25.64
CENTER	34.09	9.73	28.20
SCALE	38.33	25.33	35.55

Table 3: Comparison of F1-scores on the multi-cause scenario.

	RECCON-DD	RECCON-IE	ECF
Full model	58.83	34.69	57.70
w/o SRL	56.72	31.22	55.77
w/o GA	55.26	29.08	53.82
w/o SRL & GA	53.66	28.18	52.54
w/o EE	58.44	34.43	57.43
w/o CE	57.99	34.07	57.11
w/o CE & EE	57.15	33.57	56.54

Table 4: Ablation study.

multi-cause prediction. Nevertheless, SCALE consistently achieves the best F1 scores across all three subsets. We attribute these improvements to the global alignment formulation in SCALE, which models emotion-cause relations as soft many-to-many correspondences between emotion-oriented and cause-oriented representations, enabling joint reasoning over dispersed causal evidence for each emotion.

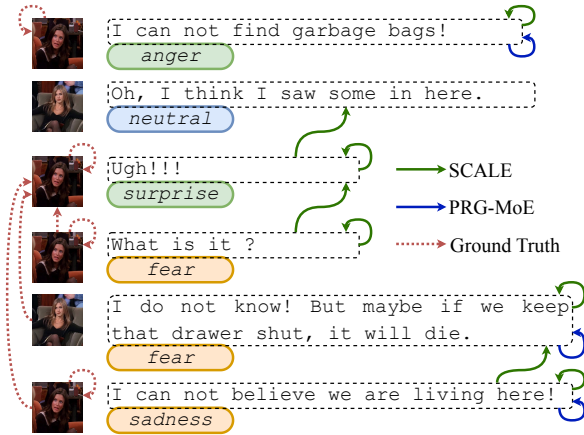
3.2.3 Ablation Study (RQ3)

To verify the effectiveness of the key design principles in SCALE, we conduct ablation experiments by removing separated representation learning (SRL; see §2.2), global alignment (GA; see §2.3), and auxiliary supervision (i.e., EE and CE). Specifically, w/o SRL collapses emotion-oriented

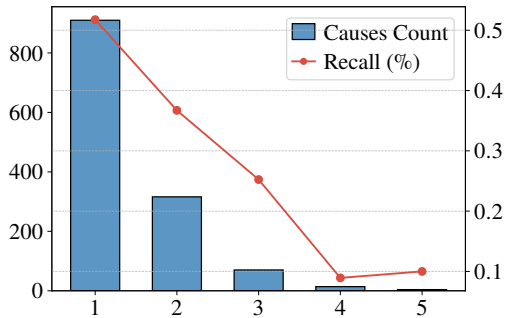
and cause-oriented encoders into a single graph encoder that learns a shared representation for all utterances, while w/o GA removes the alignment module and performs emotion-cause prediction based solely on independent pairwise scores. As shown in Table 4, removing either SRL or GA consistently degrades ECPEC F1 performance across all datasets, with a more pronounced drop observed when GA is disabled, underscoring the importance of soft many-to-many alignment for modeling complex emotion-cause relations. We further observe that auxiliary supervision is beneficial to ECPEC. Removing EE or CE individually leads to mild performance drops, while jointly removing both results in a more noticeable degradation, indicating that EE and CE provide complementary support for representation learning.

3.2.4 Qualitative Analysis (RQ4)

Case Study Figure 3a presents a qualitative comparison between SCALE and the strongest baseline PRG-MoE on a dialogue sampled from the ECF dataset. Among the seven ground-truth emotion-cause pairs, SCALE correctly predicts five, whereas PRG-MoE identifies only two, indicating a clear performance gap. For long-distance dependencies, such as (u_6, u_3) and (u_5, u_3) , both models fail to recover the correct relations, suggesting that capturing causal cues separated by large conversational gaps remains challenging. In contrast, for short-distance dependencies, including (u_1, u_1) , (u_3, u_3) , (u_4, u_4) , (u_4, u_3) , and (u_5, u_5) , SCALE achieves perfect predictions. Regarding multi-cause scenarios, SCALE successfully identifies the dual-cause emotion $[(u_4, u_3), (u_4, u_4)]$, but fails to capture both causes in $[(u_6, u_6), (u_6, u_3)]$. These observations suggest that while the global alignment mechanism enables flexible one-to-many reasoning, modeling long-range causal dependen-



(a) Qualitative comparison.



(b) Cause-number distribution and corresponding recall on ECF.

Figure 3: Qualitative analysis.

511 cies remains an open challenge.

512 **Error Analysis.** To further analyze the behavior
 513 of SCALE in multi-cause scenarios, we examine
 514 its performance on the ECF dataset. As shown in
 515 Figure 3b, most instances involve a single cause,
 516 while samples with multiple causes are increasingly
 517 scarce. Despite achieving the best overall perform-
 518 ance on the multi-cause subset of ECF, SCALE
 519 exhibits decreasing recall as the number of causes
 520 increases, since recall requires all causes associ-
 521 ated with an emotion to be correctly identified, in-
 522 dicating that exhaustive cause retrieval in complex
 523 multi-cause settings remains challenging.

524 3.2.5 Auxiliary Analysis

525 **Performance on EE and CE.** To provide ad-
 526 ditional context on the intermediate subtasks, we
 527 evaluate the performance of SCALE on emotion ex-
 528 traction (EE) and cause extraction (CE). As shown
 529 in Table 5, SCALE yields reasonable performance
 530 on both EE and CE across datasets, without relying
 531 on task-specific architectural designs. It is worth
 532 noting that SCALE is primarily optimized for the
 533 ECPEC objective, while EE and CE are incorpo-

Method	RECCON-DD		RECCON-IE		ECF	
	EE	CE	EE	CE	EE	CE
MECPE-2steps	71.30	65.81	42.52	44.49	79.10	70.13
PRG-MoE	73.86	-	57.29	-	71.82	-
Joint-Xatt	58.71	51.93	47.34	40.26	67.75	62.73
Joint-GCN	61.87	51.35	46.58	35.75	68.31	64.05
MRC	75.49	-	39.54	-	74.08	-
CENTER	68.32	-	49.61	-	67.07	-
SCALE	73.23	67.87	55.57	54.90	76.10	61.81

Table 5: Comparative results of EE and CE subtasks (F1-score) across three datasets.

Method	RECCON-DD	ECF
DeepSeek-V3.2	47.11	42.81
GPT-5.1 Instant	55.26	54.76
Gemini-3-pro-preview	56.08	55.42
SCALE	58.83	57.70

Table 6: Comparison with recent LLMs.

534 rated as auxiliary supervision during training rather
 535 than standalone targets.

536 **Comparison with recent LLMs.** We compare
 537 SCALE with three recent LLMs under a 4-shot
 538 prompting setting. As shown in Table 6, SCALE
 539 achieves higher F1-scores than the best-performing
 540 LLM on both datasets. These results suggest that
 541 explicit modeling of conversational structure and
 542 emotion–cause relations remains advantageous for
 543 ECPEC, even in the presence of strong prompt-
 544 based LLM baselines.

545 4 Conclusion

546 In this paper, we proposed **SCALE**, a semantic
 547 alignment framework for emotion–cause pair ex-
 548 traction in conversations. By decoupling emotion-
 549 oriented and cause-oriented semantics and mod-
 550 eling their interactions through global alignment,
 551 SCALE reformulates ECPEC as a many-to-many
 552 reasoning problem over conversational structure.
 553 This design enables more holistic modeling of com-
 554 plex causal dependencies beyond independent pair-
 555 wise prediction. Extensive experiments on three
 556 benchmark datasets demonstrate that SCALE con-
 557 sistently outperforms existing approaches, particu-
 558 larly in challenging multi-cause scenarios, validat-
 559 ing its effectiveness and robustness.

560 Limitations

561 The proposed **SCALE** focuses on textual conver-
 562 sations and does not incorporate other modalities

673	Dubrovnik, Croatia. Association for Computational Linguistics.	
674		
675	Sophia Yat Mei Lee, Ying Chen, and Chu-Ren Huang.	
676	2010. A text-driven rule-based system for emotion cause detection . In <i>Proceedings of the NAACL HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text</i> , pages 45–53, Los Angeles, CA. Association for Computational Linguistics.	
677		
678		
679		
680		
681		
682	Bobo Li, Hao Fei, Fei Li, Tat-seng Chua, and Donghong Ji.	
683	2024. Multimodal emotion-cause pair extraction with holistic interaction and label constraint . <i>ACM Trans. Multimedia Comput. Commun. Appl.</i>	
684		
685		
686	Raymond Li, Samira Ebrahimi Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal.	
687	2018. Towards deep conversational recommendations. In <i>Advances in Neural Information Processing Systems</i> , volume 31. Curran Associates, Inc.	
688		
689		
690		
691	Wei Li, Yang Li, Vlad Pandelea, Mengshi Ge, Luyao Zhu, and Erik Cambria.	
692	2023. Ecpec: Emotion-cause pair extraction in conversations . <i>IEEE Transactions on Affective Computing</i> , 14(3):1754–1765.	
693		
694		
695	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu.	
696	2017. Dailydialog: A manually labelled multi-turn dialogue dataset. In <i>IJCNLP 2017</i> .	
697		
698		
699	Tingting Liang, Chenxin Jin, Lingzhi Wang, Wenqi Fan, Congying Xia, Kai Chen, and Yuyu Yin.	
700	2024. Llm-redial: A large-scale dataset for conversational recommender systems created from user behaviors with llms . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 8926–8939, Bangkok, Thailand. Association for Computational Linguistics.	
701		
702		
703		
704		
705		
706		
707	Bing Liu.	
708	2022. <i>Sentiment Analysis and Opinion Mining</i> . Springer Nature.	
709	Chen Liu, Changyong Niu, Jinge Xie, Yuxiang Jia, and Hongying Zan.	
710	2023. Emotion-cause pair extraction in conversations based on multi-turn mrc with position-aware gcn . In <i>2023 International Conference on Asian Language Processing (IALP)</i> , pages 25–30.	
711		
712		
713		
714		
715	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov.	
716	2019. Roberta: A robustly optimized bert pretraining approach . Preprint, arXiv:1907.11692.	
717		
718		
719		
720	Ilya Loshchilov and Frank Hutter.	
721	2017. Decoupled weight decay regularization. <i>arXiv preprint arXiv:1711.05101</i> .	
722		
723	Navonil Majumder, Soujanya Poria, Devamanyu Hazarika, Rada Mihalcea, Alexander Gelbukh, and Erik Cambria.	
724	2019. Dialoguernn: An attentive rnn for emotion detection in conversations . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 6818–6825.	
725		
726		
727		
728		
	Hermina Petric Maretic and Mireille EL Gheche.	
	2019. Got: An optimal transport framework for graph comparison . In <i>33rd Conference on Neural Information Processing Systems (NeurIPS 2019)</i> .	
	Cam-Van Thi Nguyen, The-Son Le, Anh-Tuan Mai, and Duc-Trong Le.	
	2024. Ada2i: Enhancing modality balance for multimodal conversational emotion recognition . In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 9330–9339, Melbourne VIC Australia. ACM.	
	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others.	
	2019. <i>PyTorch: an imperative style, high-performance deep learning library</i> . Curran Associates Inc., Red Hook, NY, USA.	
	Maria Pontiki, Dimitris Galanis, Haris Papageorgiou, Ion Androutsopoulos, Suresh Manandhar, Mohammad AL-Smadi, Mahmoud Al-Ayyoub, Yanyan Zhao, Bing Qin, Orphée De Clercq, Véronique Hoste, Marianna Apidianaki, Xavier Tannier, Natalia Loukachevitch, Evgeniy Kotelnikov, Nuria Bel, Salud María Jiménez-Zafra, and Gülşen Eryiğit.	
	2016. Semeval-2016 task 5: Aspect based sentiment analysis . In <i>Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)</i> , pages 19–30, San Diego, California. Association for Computational Linguistics.	
	Soujanya Poria, Navonil Majumder, Devamanyu Hazarika, Deepanway Ghosal, Rishabh Bhardwaj, Samson Yu Bai Jian, Pengfei Hong, Romila Ghosh, Abhinaba Roy, Niyati Chhaya, Alexander Gelbukh, and Rada Mihalcea.	
	2021. Recognizing emotion cause in conversations . <i>Cognitive Computation</i> , 13(5):1317–1332.	
	Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau.	
	2019. Towards empathetic open-domain conversation models: A new benchmark and dataset . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5370–5381, Florence, Italy. Association for Computational Linguistics.	
	Shruti Saxena and Joydeep Chandra.	
	2024. A survey on network alignment: Approaches, applications and future directions . In <i>Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence</i> , pages 8216–8224, Jeju, South Korea. International Joint Conferences on Artificial Intelligence Organization.	
	Konstantinos Skitsas, Karol Orłowski, Judith Hermanns, Davide Mottin, and Panagiotis Karras.	
	2023. Comprehensive evaluation of unrestricted graph alignment algorithms .	
	Huynh Thanh Trung, Nguyen Thanh Toan, Tong Van Vinh, Hoang Thanh Dat, Duong Chi Thang, Nguyen	

787	Quoc Viet Hung, and Abdul Sattar. 2020. A comparative study on network alignment techniques . <i>Expert Systems with Applications</i> , 140:112883.	emotion-cause pair extraction . <i>Knowledge-Based Systems</i> , 286:111342.	843 844
790	Botao Wang, Keke Tang, and Peican Zhu. 2024a. Enhancing emotion-cause pair extraction in conversations via center event detection and reasoning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10773–10783, Miami, Florida, USA. Association for Computational Linguistics.	A Related Work	845
791		A.1 Emotion-Cause Pair Extraction.	846
792		Research on emotion-cause analysis originated	847
793		from the Emotion Cause Extraction (ECE)	848
794		task (Lee et al., 2010; Gui et al., 2016; Li et al.,	849
795		2018), which aims to identify the cause span cor-	850
796		responding to a given emotion. To overcome the	851
797	Fanfan Wang, Zixiang Ding, Rui Xia, Zhaoyu Li, and	limitation of requiring emotion annotation before	852
798	Jianfei Yu. 2023a. Multimodal emotion-cause pair	cause extraction, Xia and Ding (2019) proposed	853
799	extraction in conversations . <i>IEEE Transactions on</i>	the Emotion-Cause Pair Extraction (ECPE) task,	854
800	<i>Affective Computing</i> , 14(3):1832–1844.	which jointly extracts emotions and their causes,	855
801	Yan Wang, Bo Wang, Yachao Zhao, Dongming Zhao,	and inspired a line of subsequent studies (Ding	856
802	Xiaojia Jin, Jijun Zhang, Ruifang He, and Yuex-	et al., 2020; Chen et al., 2020; Fan et al., 2021).	857
803	ian Hou. 2024b. Emotion recognition in conversa-	While remarkable progress has been made in ECPE	858
804	tion via dynamic personality. In <i>Proceedings of the</i>	research, the majority of existing approaches (Fan	859
805	<i>2024 Joint International Conference on Computa-</i>	et al., 2021; Chen et al., 2022; Zhu et al., 2024)	860
806	<i>tional Linguistics, Language Resources and Eval-</i>	are confined to document-level corpora, where	861
807	<i>uation (LREC-COLING 2024)</i> , pages 5711–5722,	emotions and their corresponding causes are ex-	862
808	Torino, Italia. ELRA and ICCL.	pressed within a single, coherent narrative flow.	863
809	Yejiang Wang, Yuhai Zhao, Zhengkui Wang, and Ling	Such clause-level formulations inherently neglect	864
810	Li. 2023b. Galopa: Graph transport learning with op-	the distinctive characteristics of dialogues, such	865
811	timal plan alignment. In <i>37th Conference on Neural</i>	as speaker role alternation, intertwined emotional	866
812	<i>Information Processing Systems (NeurIPS 2023)</i> .	events, and long-range conversational dependen-	867
813	Rui Xia and Zixiang Ding. 2019. Emotion-cause pair	cies, thereby highlighting the necessity of extend-	868
814	extraction: A new task to emotion analysis in texts .	ing ECPE into conversational contexts.	869
815	In <i>Proceedings of the 57th Annual Meeting of the As-</i>	A.2 Emotion-Cause Pair Extraction in	870
816	<i>sociation for Computational Linguistics</i> , pages 1003–	Conversation.	871
817	1012, Florence, Italy. Association for Computational	Early efforts on conversational cause analysis be-	872
818	Linguistics.	gan with RECCON (Poria et al., 2021), which fo-	873
819	Hongteng Xu, Dixin Luo, and Lawrence Carin. 2019.	cused on cause recognition rather than emotion-	874
820	Scalable gromov-wasserstein learning for graph par-	cause pair extraction. Li et al. (2023) formally	875
821	tioning and matching. In <i>Advances in Neural In-</i>	introduced the ECPEC task and released the Con-	876
822	<i>formation Processing Systems</i> , volume 32. Curran	vECPE dataset, along with a two-step framework	877
823	Associates, Inc.	that explicitly models conversational properties	878
824	Zhichen Zeng, Boxin Du, Si Zhang, Yinglong Xia, Zhin-	such as context dependence and speaker interactiv-	879
825	ing Liu, and Hanghang Tong. 2024. Hierarchical	ity. Subsequent works explored more sophisticated	880
826	multi-marginal optimal transport for network align-	modeling, such as pair-relations-guided mixture-of-	881
827	ment . <i>Proceedings of the AAAI Conference on Artifi-</i>	experts system PRG-MOE (Jeong and Bak, 2023),	882
828	<i>cial Intelligence</i> , 38(15):16660–16668.	machine reading comprehension-based method	883
829	Zhichen Zeng, Si Zhang, Yinglong Xia, and Hanghang	MRC (Liu et al., 2023), global-view speaker-aware	884
830	Tong. 2023. Parrot: Position-aware regularized opti-	frameworks GSESE (An et al., 2023), and event-	885
831	mal transport for network alignment . In <i>Proceedings</i>	guided ECPEC frameworks CENTER (Wang et al.,	886
832	<i>of the ACM Web Conference 2023</i> , pages 372–382,	2024a). Beyond text, multimodal ECPEC has been	887
833	Austin TX USA. ACM.	studied with the ECF dataset (Wang et al., 2023a)	888
834	Peixiang Zhong, Chen Zhang, Hao Wang, Yong Liu, and	and improved by cross-modality interaction mech-	889
835	Chunyan Miao. 2020. Towards persona-based empa-	anisms such as HiLo (Li et al., 2024). Despite	890
836	thetic conversational models . In <i>Proceedings of the</i>	these advances, existing methods still follow the	891
837	<i>2020 Conference on Empirical Methods in Natural</i>		
838	<i>Language Processing (EMNLP)</i> , pages 6556–6566,		
839	Online. Association for Computational Linguistics.		
840	Peican Zhu, Botao Wang, Keke Tang, Haifeng		
841	Zhang, Xiaodong Cui, and Zhen Wang. 2024.		
842	A knowledge-guided graph attention network for		

pairwise classification paradigm and thus fail to capture global many-to-many alignments between emotions and causes in dialogues, leaving robust causal modeling an open challenge.

A.3 Graph Alignment and Optimal Transport.

Graph alignment aims to identify correspondences between nodes across related graphs (Saxena and Chandra, 2024; Skitsas et al., 2023; Trung et al., 2020), a problem widely studied in network analysis and data integration. As a global matching problem, it can be naturally addressed by Optimal Transport (OT) (Zeng et al., 2024; Wang et al., 2023b; Xu et al., 2019), which computes a global coupling between two sets that minimizes transportation cost and captures many-to-many correspondences under a global optimization objective (Zeng et al., 2023; Margetic and Gheche, 2019). This property makes OT particularly suitable for aligning emotions and causes in dialogues, where multiple candidates may coexist and local decisions can be insufficient. However, OT has not yet been explored in ECPEC, which motivates reformulating our task as a global graph alignment problem through the integration of dual graph learning and OT.

B Datasets

We employ three representative datasets, followed by detailed descriptions below.

RECCON (Poria et al., 2021) is a widely used benchmark that consists of two subsets: **RECCON-DD**, annotated from DailyDialog (Li et al., 2017), serves as the main corpus for model training and evaluation, while **RECCON-IE**, annotated from IEMOCAP (Busso et al., 2008), is a smaller subset used exclusively to test the generalization ability.

ECF (Wang et al., 2023a) is a multimodal benchmark derived from the sitcom *Friends*, providing annotated emotion–cause pairs across text, audio, and visual modalities, with many emotions triggered by multiple utterances.

C Baselines

We compare our method against seven representative baselines, which can be broadly grouped into three categories:

1) *Method based on general pre-trained model*: Following Poria et al. (2021), we adopt a pretrained RoBERTa (Liu et al., 2019) model with a classification layer to identify emotion–cause pairs, which

Text	Audio	Video	F1
+	-	-	57.70
+	+	-	58.13
+	-	+	58.07
+	+	+	58.63

Table 7: Multimodal evaluation of SCALE on the ECF dataset.

we refer to as **RECCON**, serving as a benchmark baseline.

2) *Methods with sequential modeling*: **MECPE-2steps** (Wang et al., 2023a) adopts a two-step pipeline that first extracts candidate emotions and causes sets by a shared BiLSTM and then filters valid pairs with another BiLSTM. **PRG-MoE** (Jeong and Bak, 2023) constructs relational graph with a mixture-of-experts, where a gating network aggregates diverse relational patterns. **Joint-Xatt** (Li et al., 2023) utilize cross-attention to model emotion–cause dependencies.

3) *Methods based on graph modeling*: **Joint-GCN** (Li et al., 2023) extends Joint-Xatt by replacing cross-attention with graph convolutional network to model inter-utterance relations. **MRC** (Liu et al., 2023) reformulates ECPEC as machine reading comprehension and employs GNNs to encode dialogue structure. **CENTER** (Wang et al., 2024a) builds a center event–aware graph with contrastive objectives for pair-level discrimination.

D Multimodal Extension

Although SCALE is not primarily designed for multimodal modeling, it can be easily extended to handle multimodal inputs. As the ECF dataset provides text, audio, and video modalities, we follow prior work by concatenating unimodal features as a simple fusion strategy, since multimodal fusion is not the main focus of this paper. As shown in Table 7, SCALE achieves consistent improvements over the text-only variant when additional modalities are incorporated, demonstrating its adaptability and robustness across multimodal settings.