Multi-Party Empathetic Dialogue Generation: A New Task for Dialog Systems

Anonymous ACL submission

Abstract

Empathetic dialogue assembles emotion understanding, feeling projection, and appropri-003 ate response generation. Existing work for empathetic dialogue generation concentrates on 004 the two-party conversation scenario. Multiparty dialogues, however, are pervasive in re-007 ality. Furthermore, emotion and sensibility are typically confused; a refined empathy analysis is needed for comprehending fragile and nuanced human feelings. We address these issues by proposing a novel task called Multi-Party Empathetic Dialogue Generation in this study. A new dataset MPED with 130k multi-party dialogues is correspondingly presented for this 014 task, which makes up for the absence of a large-scale benchmark in this field. Addition-017 ally, a Static-Dynamic model for Multi-Party Empathetic Dialogue Generation, SDMPED, is introduced as a baseline by exploring the static sensibility and dynamic emotion for the multi-party empathetic dialogue learning, the aspects that help SDMPED achieve the stateof-the-art performance on MPED.

1 Introduction

024

034

040

Empathetic conversation studies have been coming to the forefront in recent years owing to the increasing interest in dialogue systems. Empathetic dialogues not only provide dialogue partners with highly relevant contents but also project their feelings and convey a special emotion, that is, empathy. As revealed by previous studies (Fraser et al., 2018; Zhou et al., 2020), empathy can enhance conversation quality and transmit appropriate emotional responses to partners. Accordingly, most, if not all, existing work focuses on taking an emotional perspective in dialogue studies (Levinson et al., 2000; Kim et al., 2004; Bertero et al., 2016; Fraser et al., 2018; Rashkin et al., 2019).

Although the empathetic conversation has received extensive attention, its exploration is still limited to the scenario with only two parties. In



Figure 1: An empathetic dialogue example of multiparty. When people with different sensibilities respond to the same requests for help, their emotions and empathy differ. Different shades of red and blue denote the degree of positive and negative emotions, and different shades of green denote the degree of sensibilities. The texts use three kinds of underlines: straight, wavy, and dotted, which depict appropriate Emotional Reactions, Interpretations, and Explorations (three criteria to assess empathy), respectively.

fact, multi-party chatting scenes are common in seminar discussions, conferences, and group chats. Multi-party conversations also rely on aid from empathy analysis. For instance, people with a similar experience can smoothly communicate with each other and easily feel understood, encouraged, and supported at a mental health support platform. These observations encourage us to present a novel natural language processing task called Multi-Party Empathetic Dialogue Generation.

Generating multi-party empathetic dialogues faces two challenges. One challenge is the way to model multi-party dialogues. First, existing two-party dialogue models follow a seq2seq structure, whereas most multi-party dialogues are nonsequential. As shown in Figure 1, in response to *Speaker 1*, the third and fourth utterances both express empathy for her stress and struggle. Second, in addition to the target participant, other participants also have implicit influence and interaction, and should be considered of generating utterances at each step. For instance, as an example of how to successfully resolve the situation, *Speaker 4* inspires *Speaker 1* as well as relieves *Speaker 3* of her worry.

061

062

065

067

069

071

073

074

077

084

086

096

100

102

103

104

105

106

107

108

109

110

Another challenge is the way to model the fragile and nuanced feelings of dialogue participants. We first clarify the relations of sensibility, emotion, and empathy in this study. Previous empathy studies recognized the emotion of one party and generated dialogues coupled with the same emotion (Rashkin et al., 2019; Shin et al., 2020). However, empathy is also determined by sensibility, which is a perspective-taking ability to experience other partners' emotions and make an appropriate response with his/her own view. According to the response "I went through this, too" in Figure 1, we can find that Speaker 4 has a similar experience to Speaker 1, while Speaker 2 can only provide superficial comfort to Speaker 1 due to his weak sensibility. We observe that sensibility arises from personality and experience, and remains static throughout a conversation. On the other hand, emotion may dynamically change. For example, Speakers 2, 3, and 4 possess different sensibilities to Speaker 1, and these personal background-related attributes are persistent in the conversation. By contrast, the emotion of Speaker 1 gets reversed after receiving positive replies, as well as the main tone of this dialogue.

We comprehensively cope with the aforementioned challenges in this study. First, we introduce a new Multi-Party Empathetic Dialogue (MPED) dataset, which contains 130k multi-label multiparty empathetic dialogues. To the best of our knowledge, MPED is the largest empathetic dialogue dataset created to date (Rashkin et al., 2019; Poria et al., 2019; Firdaus et al., 2020). Moreover, MPED covers a large number of different emotions in a balanced manner.

Furthermore, we present a Static-Dynamic model for Multi-Party Empathetic Dialogue Generation called SDMPED. SDMPED models multiparty dialogues by constructing a dynamic graph network with temporal information and explores participants' dynamic emotions and static sensibilities by fusing speaker information.

The contributions of our work are as follows:

• We propose a new task called Multi-party Em-

pathetic Dialogue Generation, which attempts to resolve the emotional changes and empathy generation of multiple participants in a conversation. We also introduce a novel large-scale multi-party empathetic dialogue dataset which paves the way for future emotion-centered dialogue examinations.

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

- We propose an effective baseline model SDMPED for this new task, which combines dynamic emotions and static sensibilities from multiple parties.
- We demonstrate that our approach leads to performance exceeding the state of the art when trained and evaluated on MPED.

2 Related Work

2.1 Empathy Analysis

Considering empathy in modeled conversations has been proposed as early as 20 years ago (Levinson et al., 2000). However, this idea has not been widely studied in NLP field due to the limitations of the available data. Recently, Rashkin et al. (2019) re-introduced the concept of empathetic dialogue and constructed the first empathetic dialogue dataset, EMPATHETICDIALOGUES (ED), which contains 32 emotions in 25K dialogues. Another dataset, PEC (Zhong et al., 2020), provides assurance that most of the data are in line with the characteristics of empathy, yet it lacks emotion-related annotations. Another limitation is that data in PEC come from only two forums on Reddit (i.e., happy5 and offmychest). The data in BlendedSkillTalk dataset (Smith et al., 2020) are collected from the ED, ConvAI2 (Dinan et al., 2020), and Persona-Chat (Zhang et al., 2018) datasets. However, only a small portion of these data are characterized by empathy. Notably, none of the aforementioned datasets have multiple (>2) persons participating in the same conversation, neither they include empathy degree labels.

Shin et al. (2020) formulated a reinforcement learning problem to maximize the user's emotional perception of the generated responses. Li et al. (2020b) utilized the coarse-grained dialogue-level and the fine-grained token-level emotions, which helped better capture the nuances of user emotions. In Caire (Lin et al., 2020), the empathy generation tasks are reinforced with an auxiliary objective for emotion classification by using a transfer learning model. Nevertheless, current empathetic dialogue



(a) Emotion Distribution (b) Empathy Distribution Figure 2: Distribution of emotions and empathy degrees on MPED.

models are conducted in the context of two participants; they do not explore the implicit interactions among multiple speaking persons and do not consider the differences in their sensibilities.

2.2 Multi-Party Dialogue

160

161

162

163

164

165

166

168

170

171

172

173

174

175

176

178

179

180

182

183

186

187

190

191

192

193

194

195

196

197

198

199

Over the last years, researchers have gradually shifted from studying simple emotions in twoparty dialogues (Busso et al., 2008; Li et al., 2017) to conducting more complex emotion analysis of multiple participants. STAC (Asher et al., 2016) and ARS (Ouchi and Tsuboi, 2016) are the multiparty dialogue datasets without emotion labels. MELD (Poria et al., 2019) and MESID (Firdaus et al., 2020) create the multi-modal multi-party emotional dialogue datasets from the TV series Friends. However, these two datasets contain the emotion-related data derived from short and colloquial chats from TV series, and consequently, their dialogue quality cannot be guaranteed. Additionally, these datasets can only be utilized for simple upstream tasks, such as emotion recognition. Most of the dialogues in current datasets are daily conversations on trivial topics, while those modeling empathy dialogues are lacking.

Majumder et al. (2019) proposed a conversational emotion recognition model based on RNN to dynamically model the states of multiple speakers. Later, Ghosal et al. (2019) and Li et al. (2020a) also studied context and speaker sensitivity based on the approach of Majumder et al. (2019). A common problem of these models is that they only focus on the accuracy of emotion recognition while ignoring the dynamic changes of emotions.

3 MPED: Multi-Party Empathetic Dialogue Dataset

In this section, we introduce the creation process of MPED and its statistics. We regard an empathetic post and its meaningful replies as a dialogue and ensure that each dialogue has more than 3 participating speakers. Our dataset includes posts that contain replies from multiple people, along with associated emotion and empathy degree labels. The empathy degree label of each utterance will be used in conjunction with the emotional content in our future model to learn the sensibility of each person. 200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

229

230

231

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

We propose a concept called dialogue emotional turn, which is different from the traditional dialogue turn. We assume that a dialogue can have multiple sentences in one emotional turn, but with the same emotional tone. When a person utters a second sentence, the emotion may already differ from the previous one. Other people's subsequent utterances and emotions will be centered around this sentence. Therefore, we divide the dialogues to study the emotion variations over time, according to the principle that the same speaker can make at most one utterance during each emotional turn.

Data Collection and Pre-Processing TalkLife¹ (talklife.co), which is the largest online peer-topeer mental health support platform, provides the data. Users on TalkLife can express their anxiety, depression, and other psychological issues (e.g., eating disorders) by chatting with experts and others who have similar experiences.

Generally, we permit the words of each utterance to range between 3 and 100, excluding emojis, which are stored separately². We discard artificially repeated characters, correct spelling errors, and standardize network language. Developing a dialogue model in high-risk environments such as mental health requires more ethical considerations (Sharma et al., 2021). Therefore, we focus our analysis on help-seeking or emotional comfortseeking conversations. As a result, the conversations with sensitive contents (e.g., serious diseases and suicide) are filtered out. In the end, we further ensure that no private information is contained in our dataset.

It is quite beneficial that emotional category labels are available in TalkLife, which saves a lot of manual work. We have confirmed their accuracy and constructed the MPED dataset with 60 kinds of emotions. We further classify these emotions for simplicity into 10 types, as shown in Figure 2. MPED includes single-turn and multi-turn dialogue data, called MPED-S and MPED-M. We randomly split them into 80% training set, 10% validation set, and 10% testing set, respectively.

¹https://www.talklife.com.

²Emotional utterances have been incorporated in MPED yet not in our proposed baseline since we focus on unimodal text in this study.

Datasets	MPED-S			MPED-M			
Statistics	train	val	test	train	val	test	
Number of dialogues	110,000	10,267	10,267	1,800	212	200	
Number of utterances	451,465	42,622	42,142	10,227	1,265	1,154	
Size of vocabulary via spaCy	71,809	18,436	18,577	8,960	3,135	2,708	
Number of speakers (Avg. (Max.))	4.02 (41)	4.01 (42)	4.06 (29)	4.33 (33)	4.42 (30)	4.34 (26)	
Length of dialogues (Avg. (Max.))	4.19 (45)	4.15 (44)	4.10 (30)	5.68 (40)	5.97 (37)	5.77 (28)	
Length of utterances (Avg. (Max.))	20.79 (100)	21.00 (96)	21.01 (94)	21.89 (99)	23.41 (97)	21.65 (100)	
Number of turns (Avg. (Max.))	-	-	-	2.09 (13)	2.19 (7)	2.12 (4)	

Table 1: Statistics for MPED-S and MPED-M. Avg.and Max. are abbreviations of Average and Maximum.

Empathetic Pre-Processing Given that empathy 248 is a complex feeling, gathering empathetic data is 249 challenging. We first remove the conversations that 250 do not contain empathetic posts, such as games and poetry. Then, we design a three-point scale 253 (0 to 2) and evaluate empathy on the basis of the 254 standard proposed by Sharma et al. (2020), where three criteria are used: Emotional Reactions (ex-255 pressing warmth and compassion), Interpretation (articulating understanding of feelings and experiences), and Exploration (exploring feelings and 259 experiences not stated in the post). Considering the large data size, manually screening dialogues is infeasible; thus, we utilize the model proposed 261 by Sharma et al. (2021) to filter out simple replies 262 and label single-turn dialogues. 263

264

265

267

270 271

272

273

276

277

278

Data Analysis and Comparison The summary of MPED is presented in Table 1. More than four speakers are usually available per dialogue on MPED to ensure sufficient participants in conversations. In terms of dialogue quality, the average utterance length of MPED is nearly 21, which is comparable to ED (Rashkin et al., 2019) but is much larger than MELD (Poria et al., 2019) whose average length is 8. MPED can be regarded as a large-scale and high-quality dataset containing multi-party dialogues. For example, only 1,000 dialogues (including a significant portion of singleturn conversations) are provided in MELD, which is generally insufficient for training deep learning models.

As shown in Figure 2, MPED incorporates ten 279 types of emotions and three degrees of empathy. The emotion types on MPED are more fine-grained 281 than that in MELD where nearly 47% emotions are Natural. People are unlikely to receive empathetic 283 remarks from others who have Natural emotions. Meanwhile, ten types are of a moderate scale. Compared with ED (Rashkin et al., 2019) whose emotions are classified into 64 types, this scale is more 287 conducive to emotion analysis. For example, the 288 most often used words corresponding to Afraid and *Terrified* are slightly different in ED. In multi-party

empathetic dialogue generation, the sensibilities of different people are comprehensively analyzed on the basis of emotion and empathy labels. This could not be achieved in prior work. Moreover, the relatively small proportion of *Strong degree of Empathy* in Figure 2 (b) illustrates that empathetic dialogues that can be felt and truly empathized are relatively rare, as it is in reality. 291

292

293

294

295

296

297

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

4 Model

In this section, we introduce a static-dynamic model called SDMPED as shown in Figure 3. We begin by describing the construction of the Temporal Dynamic Graph Network (TDGCN), including speaker sensibility nodes, emotion-related utterance nodes, and various types of edges between them. Thereafter, we use TDGCN to obtain dynamic emotions and static speaker sensibilities by integrating nodes and edges. Finally, we use prompt tuning to generate final dialogue responses based on emotion and sensibility information.

4.1 **Problem Definition**

First, we introduce key symbols and concepts used in our study. A T emotional turns dialogue with Nutterances between M (M > 2) speakers can be expressed as $U = \{u_{ik} | 1 \le i \le N \text{ and } 1 \le k \le M\}$, where u_{ik} represents the *i*th sentence from *j*th speaker. To better study emotion variations, we specify that a speaker can at most utter one sentence in each emotional turn. Thus, U can be divided into $U = \{U_t | 1 \le t \le T\}$, where each part U_t has n_t nodes. Further, the sensibilities of speakers can be expressed as $S = \{s_1, s_2, ..., s_M\}$. Our model aims to generate an empathy response of length L.

4.2 Graph Construction

SDMPED captures the sensibility information and emotional variations of multiple parties owing to a novel graph network.

First, we train the multi-scale TextCNN (Zhang and Wallace, 2015) according to the empathy



Figure 3: The overall architecture of SDMPED. Feature extraction provides the utterance and speaker sensibility nodes u_j and s_i , which will be input into TDGCN. By considering the utterance nodes and a segmented edge matrix E_t at time step t, we are able to compute the emotion-related content features. We combine static sensibilities with the current content information to get dynamic emotional information and input into the next moment. Finally, we use prompt tuning to generate final dialogue responses based on the dynamic emotions at t + 1.

degrees of our dataset, and we extract the *d*dimensional utterance-level features containing sensibility information. In each turn, we use the emotion of the first speaker as the main emotional tone, and extract the emotional content features based on those emotion labels in the same way.

Using these sensibility-related features as nodes and speaker-utterance relationships as an adjacency matrix, we construct a two-step static graph network to determine the static sensibility information $H_S = \{(H_x)_S | 1 \le x \le M\}$ of speakers. Thereafter, we represent the dialogue as a directed graph G = (V, E, R) to obtain additional emotional information. The graph is constructed as follows:

Nodes V: The node set $V = \{v_{ik} | 1 \le i \le N \text{ and } 1 \le k \le M\}$ incorporates emotion-related utterances. Among them, each node v_{ik} (abbreviated as v_i) is initialized with the extracted feature u_i spoken by the speaker s_k .

Adjacency Matrix E: $e_{ij} \in E$ represents the edge from the utterance node v_i to v_j . Before feeding it into TDGCN, we need to divide E into T steps: $E = \{E_t | 1 \le t \le T\}$. At time step t, the divided matrix E_t includes only edges corresponding to the utterance in the emotional turn t.

Edge Relations R: The relationship r_{ij} of edge e_{ij} is set mainly depending upon two things (Ghosal et al., 2019; Yang et al., 2021): the relative occurrence positions of u_i and u_j in the conversation (with three types of relations, namely, *Before*, *Current*, and *After*) and both speakers of the constituting utterance nodes, as shown in Figure 4.

As shown in Figure 1, four speakers participate

in the dialogue with 7 utterances. This dialogue has two emotional turns: u_1 to u_4 and u_5 to u_7 . The nodes and edges are constructed in Figure 4. We take node u_3 as an example. The edge e_{13} represents that u_1 spoken by s_1 appears before u_3 spoken by s_3 and the influence between them; the self-loop e_{33} represents the influence of current node u_3 on itself. 364

365

366

367

368

369

370

371

372

373

374

375

377

378

379

380

381

382

383

388

389

390

391

392

Two-Step Graph Update: Utilizing the Two-Step Graph Update mechanism, we can effectively normalize the local neighborhood through neighborhood connections and enable self-dependent feature transformation through self-connections, thereby extracting further information (Ghosal et al., 2019):

$$h_i^{(1)} = \sigma(\sum_{r \in R} \sum_{j \in N_i^r} \frac{\alpha_{ij}}{c_{i,r}} W_r^{(1)} u_j + \alpha_{ii} W_0^{(1)} u_i),$$

$$h_i^{(2)} = \sigma(\sum_{j \in N_i^r} W^{(2)} h_j^{(1)} + W_0^{(2)} h_i^{(1)}), \quad (1)$$

where α_{ij} and α_{ii} are the edge weights and N_i^r denotes the neighboring indices of node v_i under relation $r \in R$ and $c_{i,r} = |N_i^r|$. σ is the activation function ReLU, while $W_r^{(1)}$, $W_0^{(1)}$, $W^{(2)}$, and $W_0^{(2)}$ are learnable parameters. We can call these two steps RGCONV and GCONV respectively.

4.3 TDGCN

Previous dynamic graphs were mostly used in spatio-temporal traffic networks with separated spatial and time features (Guo et al., 2019; Zhao et al., 2020). However, given that the utterance node is time-related and changes frequently, we implement the dynamic graph by updating a weight ma-

363



Figure 4: Transformation of dynamic emotions from t_1 to t_2 , as well as various types of edges between different nodes (e.g., Node u_3).

trix through GRU and updating the hidden layer through the two-step graph:

$$M_t^{(l)} = \text{GRU}(H_{t-1}^{(l)}, M_{t-1}^{(l)}),$$

$$H_t^{(l)} = \text{GCONV}(\text{RGCONV}(E_t, H_{t-1}^{(l)}, M_t^{(l)})), (2)$$

where $t \in [1, T]$ and $l \in [1, L]$ (*L* generally equals 2) denote the time and layer index, respectively. $M_{t-1}^{(l)}$ represents the weight matrix updated by GRU. $H_t^{(0)}$ is equal to the node features V. The hidden state $H_t^{(l)}$ of the *l*th layer at time step *t* can be divided into n_t parts: $H_t^{(l)} = \{(h_x)_t^{(l)}\}$, where *x* represents the speaker index. By concatenating person's sensibility with corresponding emotionrelated content $(h_x)_t^{(l)}$, we obtain dynamic emotion embedding:

$$(e_x)_t^{(l)} = \left[(H_x)_S; (h_x)_t^{(l)} \right].$$
 (3)

Then, the emotion embedding set $e_t = \{(e_x)_t^{(l)}\}$ is sent to a fully connected layer and regarded as H_t at t + 1 time step. We can also obtain a crossentropy loss function at t + 1:

$$P_e = softmax(W_le_{t+1}), L_{emo} = -\log\left(P_e[e]\right).(4)$$

4.4 Decoder and Loss

We adopt prompt tuning (Lester et al., 2021) to generate responses, which is a lightweight alternative to fine-tuning the generation task and keeps language model parameters unchanged while optimizing the prompt. The prompt adjustment achieves comparable performance in the full data setting by learning only parameters with a small proportion.

The representation e_{t+1} is first transformed by a linear transformation into prompt. We can obtain the input of the empathy decoder Z = [X; prompt; Y], where X and Y represent the context and target response, respectively. We use the standard maximum likelihood estimate to optimize the response prediction, and we obtain another loss function through the decoder: 423

424

425

426

427

428

429

430

431

432 433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

$$L_{res} = -\log(p(Y|R_{generate})). \quad (5)$$

Finally, all the parameters are jointly trained end-to-end to optimize the listener selection and response generation by minimizing the sum of two losses:

$$L = L_{emo} + L_{res}.$$
 (6)

5 Experiments

5.1 Experimental Setting

The hyper-parameters in our approach are set as follows. The input embeddings are 300-dimensional pre-trained 840B GloVe vectors. The speaking coefficient c is 5. The learning rate is 0.003 and batch size is 16. The dropout rate is 0.6, while the loss weight is $5e^{-4}$.

5.2 Evaluation Criteria

Automatic Evaluation Criteria We calculate the AVG BLEU (average of BLEU-1,-2,-3,-4) (Papineni et al., 2002) and ROUGE-L (Lin, 2004) scores as evaluations of model response generation, which have been often used to compare the system-generated response against the human-gold response in generation tasks.

Human Evaluation Criteria We randomly collect 100 dialogue samples and their corresponding generations from each model. Then, we assign human annotators to rate each response between 1 and 5 on three distinct attributes: *Empathy* assesses whether the speaker of the response understands the feelings of others and fully manifests it; *Relevance* evaluates whether the response is relevant with the dialogue context and topic; and *Fluency* measures whether the response is smooth and grammatically correct.

5.3 Baselines and Models

6

MReCoSa: A context-sensitive model with multihead self-attention (Zhang et al., 2019). Multi-Trans: This multi-task model learns emotion classification and dialogue generation at the same time (Rashkin et al., 2018). MoEL: This model (Lin et al., 2019) combines the response representations from multiple emotion-specific decoders. EmpGD: This method (Li et al., 2020b) exploits coarse-grained and fine-grained emotions by an adversarial learning framework. Caire: This method

395

- 39
- 399 400
- 401
- 402 403
- 404
- 405
- 406 407

408 409

410

411

412

413

414

415

416

417

418

419

420

421

Model	MPED-M					MPED-S				
Metrics	ROUGE-L	AVG BLEU	Emp.	Rel.	Flu.	ROUGE-L	AVG BLEU	Emp.	Rel.	Flu.
MReCoSa	10.31	2.58	2.20	3.09	3.91	10.74	3.90	2.22	3.34	4.00
Multi-Trans	6.59	3.86	2.81	3.13	3.92	8.10	4.22	2.76	3.41	4.20
MoEL	6.83	2.99	3.11	3.07	3.89	8.44	3.13	3.00	3.28	4.13
EmpDG	10.86	4.26	3.19	3.39	4.30	11.53	4.52	3.32	3.55	4.30
Caire	11.58	4.85	3.17	3.62	4.37	12.48	5.49	3.30	3.89	4.46
Random prompt	11.36	4.68	3.10	3.65	4.10	12.04	5.41	3.44	3.81	4.40
SDMPED w/o S	12.06	5.57	3.29	3.66	4.30	13.47	5.88	3.51	3.81	4.53
SDMPED	12.87	6.35	3.40	3.74	4.39	14.16	7.37	3.71	3.86	4.59

Table 2: Experimental results on MPED. The automatic evaluations include AVG BLEU and ROUGE-L, and Emp.; Rel. and Flu. stand for the human evaluations *Empathy*, *Relevance* and *Fluency*.

Model	MPED	-M	MPED-S		
Matrice	DOLICE I	AVG	POLICE I	AVG	
Metrics	KOUGE-L	BLEU	KOUGE-L	BLEU	
SDMPED	12.87	6.35	14.16	7.37	
SDMPED w/o S	12.06	5.57	13.17	5.88	
Two-Step Graph	11.54	4.87	12.39	5.69	
Graph-Based	11.23	4.67	11.68	4.84	

Table 3: Ablation study on MPED-M and MPED-S.



Figure 5: The effect of different numbers of speakers. The orange and blue lines represent BLEU-1 and ROUGE-L, and histograms in dark blue show the average number of words spoken by each person in multi-turn dialogues.

(Lin et al., 2020) fine-tunes a large-scale pre-trained language model with multiple objectives: response language modeling, response prediction, and dialogue emotion detection. **Random Prompt:** We built a network with random values for prompt according to Lester et al. (2021).

We describe the variants of our model below: **Graph-Based:** This simple model uses a graphbased model to build the empathetic dialogue graph of multi-party. **Two-Step Graph:** This model adopts a graph network with two-step graph update. **SDMPED without Sensibility (SDMPED w/o S):** This model ignores the sensibilities of speakers but maintains a TDGCN structure. **SDMPED:** Our final model combines dynamic emotions with static sensibilities to produce empathy responses.

5.4 Experimental Results

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

Automatic Evaluation Results According to the experimental results shown in Table 2, our model SDMPED achieves the highest scores under most metrics compared with other baselines. The noticeable improvement indicates the effectiveness of



Figure 6: The effect of different numbers of tokens. The first three lines of this legend compare the effects when the emotion categories are 6, 10, and 60. **Before Utterance** and **Before Response** compare the effects of using different prompt embedding positions when dividing emotions into 10 categories.

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

SDMPED on empathetic expressions of multi-party. Since multi-party dialogues are not time-sequential and multi-turn dialogues need to consider the impact of each turn, SDMPED performs better than the models MoEL, EmpDG, and Caire that are designed solely for two-party dialogue. Compared with the Random prompt model, our model has been greatly improved, which demonstrates that our emotional prompt design plays an important role. Given that persons have different sensibilities, adding the characteristics of different people to explore their conversations helps improve the performance. Thus, SDMPED obtains a performance improvement on the basis of SDMPED without Sensibility.

Human Evaluation Results Table 2 shows that SDMPED has achieved good performance in *Empathy, Relevance,* and *Fluency.* Our model is effective in capturing different emotional changes between multiple speakers and generating appropriate responses. MoEL and EmpDG are more inclined towards the characteristics of two-party dialogues, and thus cannot fully adapt to the new situation of multi-party. Random prompt and Caire are basically as good as our model in *Fluency*, however their *Empathy* and *Relevance* are inferior. These two models are pre-trained transfer learning models, and the generated responses are fluent and grammatical while being simple and general.

	Speaker	Sensibility	Utterance						
Context	Lily	-	Today I broke up with my boyfriend . We had a very toxic relationship and						
			I decided to end it. Now I am alone and I don't have any friends . Who can						
			give me a single hug? (Depressed)						
	Numb	Weak	A virtual, because it could be possible. (Calm)						
Response	Eldar	Moderate	Lots of hugs to you. You can see me as your friend if you need to talk. (Worried)						
	Jain	Strong	I am truly sorry that today sucks.Enjoy yourself and focus on what you love todo, I believe you will get through it.(Optimistic)						
	Calista	Strong	Making the right decision <u>rather than staying miserable!</u> You will <u>end up finding</u> one that treats you amazingly. <u>Sending you sunshine to brighten your day</u> . (Supportive)						

Table 4: An example of different responses by different speakers. Shades of blue represent the attention weights of *Calista*. Below the text are three kinds of lines: straight, wavy, and dotted, which depict appropriate Emotional Reactions, Interpretations, and Explorations (three criteria to assess empathy).

5.5 Ablation Study

523

524

525

526

527

529

530

531

533

534

535

536

537

538

539

540

541

542

544

545

546

547

548

551

552

553

554

555

556

We perform an ablation study to better understand the contributions of the main parts of our model. As shown in Table 3, the performance becomes noticeably worse, especially in the multi-turn dialogue data, after we remove the sensibility component. The degree of empathy for empathetic dialogues depends on the emotional tone at that time and the speakers' own abilities of perspective-taking, so studying sensibilities can help better investigate the responses generated by different people. According to the comparison of SDMPED without Sensibility and Two-Step Graph, emotions of people change at every moment, and updating the graph structure at each emotional turn is particularly necessary. After removing the two-step graph update mechanism, we find that the results of Graph-Based have further declined, which indicates that the two-step graph convolution process can better extract empathetic and dialogue features.

5.6 Analysis of Speakers and Tokens

We investigate the effects of different numbers of speakers and tokens. When 3–7 speakers are available, as shown in Figure 5, the model maintains fairly stable results, indicating that it can handle multiple-party empathetic dialogues effectively. However, the results decline as the speaker number continues to increase. The reason for the drop is that our conversations are typically concentrated between 3 to 5 people, and those with more than 7 people contain little content per speaker.

In Figure 6, we compare our model with two prompt embedding methods and different numbers of emotion classification categories. The comparison between the orange and blue curves shows that dividing emotions into 10 categories gives better results than the 6 and 60 categories (6 and 60 categories similar to the number of categories in MELD and ED datasets). Clearly, dividing emotions into 10 categories and placing a prompt matrix with 2 tokens before the response can yield promising performance. 559

560

561

562

563

565

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

587

589

590

591

5.7 Case Study

We apply different speakers' sensibilities to the empathy decoder in the same multi-turn conversation context and obtain results based on MPED in Table 4. When presented with Lily's loneliness and depression, the following four speakers are willing to provide support, but they come up with different responses due to their different sensibilities. Numb is relatively unable to appreciate the emotions of Lily and jokes that she can find a virtual friend to hug; Eldar expresses warmth and suggests Lily can consider herself as a friend. Jain and Calista comfort Lily and express their understanding of how she feels after breaking up with her boyfriend. They also look forward to the future by suggesting that Lily can do something she likes to distract herself and believe that she can find the right person.

6 Conclusions and Future Work

We have introduced a novel task called Multi-Party Empathetic Dialogue Generation and a large-scale dataset, i.e., MPED. We have proposed a model called SDMPED suitable for the characteristics of the task. Our experiments have demonstrated that SDMPED is superior to other approaches on MPED. Future work can explore related issues such as integrating empathy into the dialogues, combining emojis and responses, guiding the active development of conversation.

References

593

594 595

596

597

598

600

606

607

610

611

612

613

614

619

621

622

631

632

633

637

641

647

648

- Nicholas Asher, Julie Hunter, Mathieu Morey, Farah Benamara, and Stergos Afantenos. 2016. Discourse structure and dialogue acts in multiparty dialogue: the stac corpus. In *Proceedings of the 10th International Conference on Language Resources and Evaluation*, pages 2721–2727.
- Dario Bertero, Farhad Bin Siddique, Chien-Sheng Wu, Yan Wan, Ricky Ho Yin Chan, and Pascale Fung. 2016. Real-time speech emotion and sentiment recognition for interactive dialogue systems. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1042– 1047.
- Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. In *Proceedings of the 6th International Conference on Language Resources and Evaluation*, pages 335–359.
- Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, et al. 2020. The second conversational intelligence challenge (convai2). In *Proceedings of the* 33rd Conference on Neural Information Processing Systems Competition, pages 187–208.
 - Mauajama Firdaus, Hardik Chauhan, Asif Ekbal, and Pushpak Bhattacharyya. 2020. MEISD: A multimodal multi-label emotion, intensity and sentiment dialogue dataset for emotion recognition and sentiment analysis in conversations. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4441–4453.
 - Jamie Fraser, Ioannis Papaioannou, and Oliver Lemon. 2018. Spoken conversational AI in video games: Emotional dialogue management increases user engagement. In *Proceedings of the 18th International Conference on Intelligent Virtual Agents*, pages 179– 184.
 - Deepanway Ghosal, Navonil Majumder, Soujanya Poria, Niyati Chhaya, and Alexander Gelbukh. 2019.
 DialogueGCN: A graph convolutional neural network for emotion recognition in conversation. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, pages 154–164.
- Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. 2019. Attention based spatialtemporal graph convolutional networks for traffic flow forecasting. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, pages 922– 929.
 - Sung Soo Kim, Stan Kaplowitz, and Mark V Johnston. 2004. The effects of physician empathy on

patient satisfaction and compliance. *Evaluation & the Health Professions*, 27(3):237–251.

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

669

670

671

672

673

674

675

676

677

678

679

680

681

682

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

702

- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.
- Wendy Levinson, Rita Gorawara-Bhat, and Jennifer Lamb. 2000. A study of patient clues and physician responses in primary care and surgical settings. *The Journal of the American Medical Association*, 284(8):1021–1027.
- Jingye Li, Donghong Ji, Fei Li, Meishan Zhang, and Yijiang Liu. 2020a. Hitrans: A transformer-based context-and speaker-sensitive model for emotion detection in conversations. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4190–4200.
- Qintong Li, Hongshen Chen, Zhaochun Ren, Pengjie Ren, Zhaopeng Tu, and Zhumin Chen. 2020b. EmpDG: Multi-resolution interactive empathetic dialogue generation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4454–4466.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. DailyDialog: A manually labelled multi-turn dialogue dataset. In Proceedings of the 8th International Joint Conference on Natural Language Processing, pages 986–995.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Zhaojiang Lin, Andrea Madotto, Jamin Shin, Peng Xu, and Pascale Fung. 2019. MoEL: Mixture of empathetic listeners. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, pages 121– 132.
- Zhaojiang Lin, Peng Xu, Genta Indra Winata, Farhad Bin Siddique, Zihan Liu, Jamin Shin, and Pascale Fung. 2020. Caire: An end-to-end empathetic chatbot. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence*, 09, pages 13622– 13623.
- Navonil Majumder, Soujanya Poria, Devamanyu Hazarika, Rada Mihalcea, Alexander Gelbukh, and Erik Cambria. 2019. DialogueRNN: An attentive rnn for emotion detection in conversations. In *Proceedings* of the 33rd AAAI Conference on Artificial Intelligence, pages 6818–6825.
- Hiroki Ouchi and Yuta Tsuboi. 2016. Addressee and response selection for multi-party conversation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2133–2143.

782

783

784

785

786

760

- 704 705
- 70 70
- 708
- 710 711
- 7
- 713
- 7
- 716
- 718 719 720
- 777777
- 725 726 727 728 729
- 730 731 732
- 733 734 735
- 736 737 738
- 739 740 741
- 742

744 745

- 746
- 747
- 748
- 749 750
- 751 752

753 754

754 755 756

7

758 759

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BIEU: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th annual meeting of the Association for Computational Linguistics, pages 311–318.
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2019. MELD: A multimodal multi-party dataset for emotion recognition in conversations. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 527– 536.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2018. I know the feeling: Learning to converse with empathy.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic opendomain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381.
- Ashish Sharma, Inna W Lin, Adam S Miner, David C Atkins, and Tim Althoff. 2021. Towards facilitating empathic conversations in online mental health support: A reinforcement learning approach. In *Proceedings of the Web Conference 2021*, pages 194– 205.
- Ashish Sharma, Adam Miner, David Atkins, and Tim Althoff. 2020. A computational approach to understanding empathy expressed in text-based mental health support. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, pages 5263–5276.
- Jamin Shin, Peng Xu, Andrea Madotto, and Pascale Fung. 2020. Generating empathetic responses by looking ahead the user's sentiment. In *Proceedings of the 2020 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 7989–7993. IEEE.
- Eric Michael Smith, Mary Williamson, Kurt Shuster, Jason Weston, and Y-Lan Boureau. 2020. Can you put it all together: Evaluating conversational agents' ability to blend skills. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2021–2030.
- Jianing Yang, Yongxin Wang, Ruitao Yi, Yuying Zhu, Azaan Rehman, Amir Zadeh, Soujanya Poria, and Louis-Philippe Morency. 2021. Mtag: Modaltemporal attention graph for unaligned human multimodal language sequences. In *Proceedings of the* 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1009–1021.
- Hainan Zhang, Yanyan Lan, Liang Pang, Jiafeng Guo, and Xueqi Cheng. 2019. ReCoSa: Detecting the relevant contexts with self-attention for multi-turn di-

alogue generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3721–3730.

- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 2204–2213.
- Ye Zhang and Byron Wallace. 2015. A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification. *arXiv preprint arXiv:1510.03820*.
- Ling Zhao, Yujiao Song, Chao Zhang, Yu Liu, Pu Wang, Tao Lin, Min Deng, and Haifeng Li. 2020. T-GCN: A temporal graph convolutional network for traffic prediction. *IEEE Transactions on Intelligent Transportation Systems*, 21(9):3848–3858.
- Peixiang Zhong, Chen Zhang, Hao Wang, Yong Liu, and Chunyan Miao. 2020. Towards persona-based empathetic conversational models. In *Proceedings* of the 2020 Conference on Empirical Methods in Natural Language Processing, pages 6556–6566.
- Li Zhou, Jianfeng Gao, Di Li, and Heung Yeung Shum. 2020. The design and implementation of xiaoice, an empathetic social chatbot. *Computational Linguistics*, 46(1):53–93.