# EMPATHETIC RESPONSE GENERATION VIA EMOTION CAUSE TRANSITION GRAPH

*Yushan Qian[†‡♡], Bo Wang[†*], Ting-En Lin[‡], Yinhe Zheng[‡]*
*Ying Zhu[†], Dongming Zhao[†], Yuexian Hou[†], Yuchuan Wu[‡], Yongbin Li[‡*]*

[†]State Key Laboratory of Communication Content Cognition, People's Daily Online, Beijing, China
[‡]Alibaba Group, Beijing, China
`shuide.lyb@alibaba-inc.com`

## ABSTRACT

Empathetic dialogue is a human-like behavior that requires the perception of both affective factors (e.g., emotion status) and cognitive factors (e.g., cause of the emotion). Besides concerning emotion status in early work, the latest approaches study emotion causes in empathetic dialogue. These approaches focus on understanding and duplicating emotion causes in the context to show empathy for the speaker. However, instead of only repeating the contextual causes, the real empathic response often demonstrate a logical and emotion-centered transition from the causes in the context to those in the responses. In this work, we propose an emotion cause transition graph to explicitly model the natural transition of emotion causes between two adjacent turns in empathetic dialogue. With this graph, the concept words of the emotion causes in the next turn can be predicted and used by a specifically designed concept-aware decoder to generate the empathic response. Automatic and human experimental results on the benchmark dataset demonstrate that our method produces more empathetic, coherent, informative, and specific responses than existing models.

***Index Terms***— Empathetic Dialogue, Dialogue Systems, Emotion Cause, Human Interaction

## 1. INTRODUCTION

Empathetic dialogue aims to understand the human emotional status and generate appropriate responses. Previous works have demonstrated that empathetic dialogue systems can effectively improve user experience and satisfaction in various domains, such as chit-chat [1], customer service [2, 3, 4], and psychological counseling [5]. In psychology, two primary forms of empathy are affective empathy and cognitive empathy, constituting the ideal empathy [6]. Affective empathy seeks to feel the same emotions as others, and cognitive empathy seeks to stand in someone else's situation and better understand their contextual experiences related to emotions. In empathetic dialogue research, affective empathy has been well studied, including mixture of experts [7], emotion mimicry [8], and

multi-resolution user feedback [9]. Cognitive empathy has gradually attracted the attention of scholars in recent years, including the emotion cause of the context [10, 11], external knowledge [12, 13], etc.

As an important cognitive factor, the causes of the emotion status is an integral part of human sentiment analysis [14, 15]. However, the existing empathetic dialogue methods concerning emotion causes mainly focus on causes in the current dialogue context [10, 11]. These approaches aim to understand and duplicate emotion causes in the context to show empathy for the speaker. In fact, instead of only repeating contextual causes, the real empathetic responses often demonstrate a logical and emotion-centered transition from causes in the context to those in the responses. One way to augment the emotion cause transition modeling for response generation is to introduce external knowledge with commonsense knowledge graph [12, 13]. However, the transitions of emotion causes in empathetic dialogue are often emotion-centered, which are relatively sparse or absent in the commonsense knowledge graph and difficult to be effectively searched. An example is shown on the right of Figure 1. The transition from "girlfriend" to "love" and "together" is beyond the causes in the context and is difficult to be predicted only by the current context.

To address these issues, we propose a method, named **ECTG**, to guide the generation of empathetic responses with a **E**motion **C**ause **T**ransition **G**raph. As shown in Figure 1, the proposed method consists of three stages: Graph Construction, Response Concepts Prediction, and Response Generation. The emotion cause transition graph is automatically constructed on the golden empathetic dialogue corpus, which consumes much cost and is essential in improving empathetic dialogue. We first manually annotate a span-level emotion cause dataset and exploit a pre-trained model fine-tuned on this dataset to identify emotion cause spans. Since human dialogue [16, 17, 18] naturally centers on key concepts [19, 20], we extract keywords in emotion cause spans as key concepts, which are vertices of the graph. And edges in the graph represent natural transitions between emotion causes in the dialog. Then, combined with the hierarchical context encoder and the contextual concept flow retrieved from the graph, we use the Transformer with graph attention and Insertion Transformer to jointly optimize
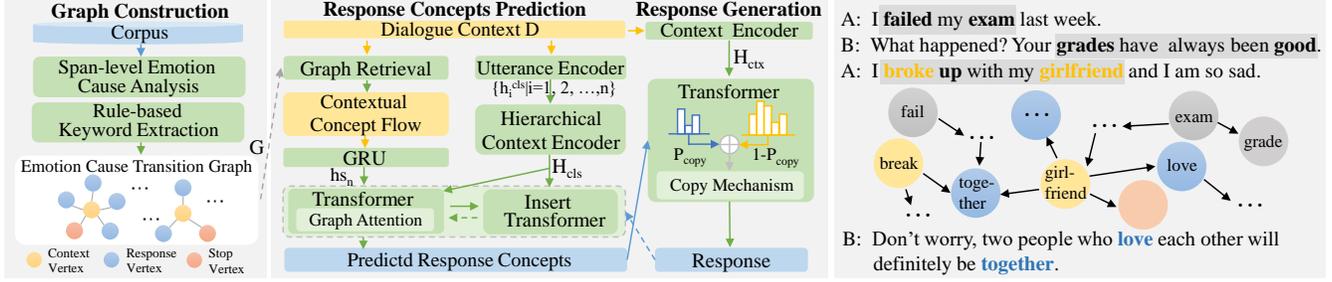
---

**Fig. 1**. The overall architecture of our proposed ECTG. The left side is the model part, and the right side is the support example.

to predict response concepts. Finally, with the dialogue context and predicted concepts, a transformer decoder with the copy mechanism explicitly generates final responses. Our contributions are summarized as follows: 1) We propose a novel approach to empathetic dialogue in line with the psychology theory and the human dialogue pattern, which can effectively improve the empathetic response generation. 2) Automatic and human evaluations show that our method generates more empathetic, coherent, informative, and specific responses than existing models. 3) To extract emotion causes more accurately, we crowdsource annotated a span-level emotion cause dataset. We will publicly release the dataset for future research.

## 2. METHODOLOGY

Formally, the constructed emotion cause transition graph is defined as $G$, given the dialog context $D$ with $n$ utterances, i.e., $D = \{U_1, U_2, \ldots, U_n\}$, $U_i$ represents the i-th utterance in $D$. Ultimately, we aim to produce empathetic, coherent, informative, and specific responses $R$.

### 2.1. Graph Construction

To construct the emotion cause transition graph, we first conduct the span-level emotion cause analysis. The emotion cause span is the consecutive sub-sequence of an utterance that expresses the cause of the emotion [21]. Due to the absence of public span-level emotion cause annotated dataset for empathetic dialogue, we follow the same setting in [10] and manually annotate the emotion cause spans in the dataset (Section 3.1).

To identify emotion cause spans, we exploit pre-trained span-level SpanBERT [22] to encode the dialog context and corresponding emotion label. We concatenate embeddings of the dialog context and emotion with the special token [SEP] as the input for the encoder. Then, we adopt Pointer network [23] to generate start and end positions of spans following [21]. We utilise the attention mechanism for each emotion cause span to measure the probability of different positions.

We identify the emotion cause span of each utterance in the dialog with the previous method. Then, we use a rule-based keyword extraction method [24] to obtain significant keywords

from emotion cause spans. All the extracted keywords are regarded as emotion cause concepts, which are defined as the vertices of the graph $G$. We connect two concepts with a direct edge if one concept appears in the last utterance of the context, which is the head vertex of the edge, and the other concept appears in the response, which is the tail vertex of the edge. We use point-wise mutual information (PMI) between the head and tail vertex to filter out low-frequency concept pairs.

### 2.2. Response Concepts Prediction

To generate empathetic responses, we predict response concepts using the emotion cause transition graph. Given the i-th utterance $U_i$, all the concepts in $U_i$ which are also involved in the graph $G$ form a concept set $cs_i = \{c_1, c_2, \cdots, c_{m_i}\}$, where $m_i$ is the number of concepts in $U_i$.

**Context Encoding.** To better utilize the dialog context [25, 26] in predicting response concepts, we encode the context hierarchically to collect all utterance representations. We prepend a special token [CLS] of each utterance $U_i$, and transform them into a sequence of hidden vectors with a BERT encoder: $h_i^{cls} = \text{BERT}_{\text{enc}}([\text{CLS}], U_i)[0]$.

$h_i^{cls}$ is the hidden representation of [CLS], which denotes the global memory of the utterance $U_i$. We input all $h_i^{cls}$ into a Transformer encoder to model the global semantic dependency between utterances: $H_{cls} = \text{Trs}_{enc}\left([h_1^{cls}, h_2^{cls}, \cdots, h_n^{cls}]\right)$.

Then, we exploit a GRU unit to recursively encode concept sets in the dialogue context:

$$hs_i = \text{GRU}\left(hs_{i-1}, \sum_{j=1}^{m_i} \alpha_{ij} e_{ij}^c\right), i \in [1, n], \quad (1)$$

$$\alpha_{ij} = \frac{\exp(\beta_{ij})}{\sum_{k=1}^{m_i}(\beta_{ik})}, \beta_{ij} = hs_{i-1}^{\mathsf{T}} W_3 e_{ij}^c, \quad (2)$$

where $e_{ij}^c$ is concept embedding, $hs_i$ represents contextual concept flow. $\alpha_{ij}$ is used to measure the probability of transitions to associated concepts.

**Response Concepts Selection.** We combine the dialogue context representation and the previously decoded concepts by a Transformer decoder, as a basis for dynamically selecting the next vertex in the emotion cause transition graph: $hdc_t =$

$\text{Trs}_{\text{dec}}\left(\left[e_{1:t-1}^{dc}\right], H_{cls}\right)$. Here, $e_{1:t-1}^{dc}$ denotes the embeddings of previously decoded concepts at step t.

For the concept set $cs_n$ of the last utterance $U_n$ in the context, we retrieve a group of subgraphs in the graph $G$, where each concept in $cs_n$ is the head vertex and its each direct neighbor vertex is the tail vertex. The subgraph $g_i = \{(c_j, c_{jk})\}_{k=1}^{N_j}, c_j \in cs_n$, where $N_j$ is the number of vertex pairs of $c_j$ in $g_i$. We employ a dynamic graph attention mechanism to calculate the subgraph vector:

$$\alpha_j = \frac{\exp(\beta_j)}{\sum_{l=1}^{m_i} \exp(\beta_l)}, \quad (3)$$

$$\beta_j = (W_4[hdc_t; hs_n])^{\mathsf{T}} \cdot \left(W_5 \sum_{k=1}^{N_j} \alpha_{jk} \left[e_j^c; e_{jk}^c\right]\right), \quad (4)$$

where $\alpha_j$ determines the choice of subgraphs. $hs_n$ incorporates information of contextual concept flow. $\alpha_{jk}$ determines which tail vertex is selected in $g_i$:

$$\alpha_{jk} = \frac{\exp(\beta_{jk})}{\sum_{l=1}^{N_j} \exp(\beta_{jl})}, \quad (5)$$

$$\beta_{jk} = (W_6[hdc_t; hs_n; e_j^c])^{\mathsf{T}} \cdot W_7 e_{jk}^c). \quad (6)$$

Finally, the chosen response concepts at step t are derived as: $P(dc_t \mid D, G, dc_{<t}) = \alpha_j \cdot \alpha_{jk}$.

**Response Concepts Refining.** From the pilot study, we found that the response concept decoder pays more attention to frequent concepts and thus lacks variety. We conjecture that supervision signals are only concept labels but the signals from the natural empathetic response should also be used simultaneously to optimize the decoder. To solve this issue, we propose an auxiliary module that takes intermediate layers of the response concept decoder as input and takes the empathetic response as output, and optimizes with the response concept decoder together via multi-task learning. In this way, the information of empathetic responses can be transported into the response concept decoder to facilitate more abundant response concept prediction. More specifically, we exploit the Insertion Transformer [27] in a non-autoregressive manner as the auxiliary loss to choose predicted concepts inspired by [20]. The loss of the Insertion Transformer is: $L_g = \frac{1}{k+1} \sum_{pos=0}^{k} \sum_{n=il}^{jl} -\log P_n^{\text{InsTrs}} \cdot w_{pos}(n)$. For more details about the Insertion Transformer, please refer to [27].

For the loss of response concepts $C$, we use negative log-likelihood loss: $L_c = \frac{1}{|C|} \sum_{t=1}^{|C|} -\log p(c_t \mid D, G, c_{<t})$.

The optimization of predicted concepts that can generate the empathetic response is determined by the weighted sum of two previous losses: $\text{Loss}_{gc} = L_g + rL_c$. Here, $r$ is the coefficient to control the impact of concept loss.

### 2.3. Empathetic Response Generation

To generate the empathetic response, we concatenate the predicted response concepts and the previous dialog context together as a sequence to the BERT encoder, and then combine a Transformer decoder with the copy mechanism to explicitly utilize it. The final generation probabilities are computed over the word vocabulary and the selected concept words:

$$H_{ctx} = \text{BERT}_{enc}(input_{D^c}), H_{dec} = \text{Trs}_{dec}(H_{ctx}), \quad (7)$$

$$P(w) = A_h \odot P_{copy} \cdot M_{src} + (1 - P_{copy})P_{gw}(w), \quad (8)$$

$$P_{copy} = \text{Sigmoid}(W_8 \cdot H_{dec}), \quad (9)$$

$$P_{gw}(w) = \text{Softmax}(W_9 \cdot H_{dec}), \quad (10)$$

where $D^c$ is the input combining the dialogue context and predicted concepts, $input_{D^c}$ is the input ids of $D^c$. $P_{copy}$ is the probability of copying a particular word from the attention distribution directly, $M_{src}$ is an indicator matrix mapping each source word to the additional vocab containing it. We apply the cross-entropy loss for training.

## 3. EXPERIMENTS

### 3.1. Experimental Setup

**Datasets & Evaluation Metrics.** We conduct experiments on the EmpatheticDialogues [28], which is a large-scale English multi-turn empathetic dialogue benchmark dataset. For automatic metrics, we adopt BLEU-4 (B-4), BERTscore F1 ($F_{\text{BERT}}$) [29], Distinct-n (Dist-1/2), ROUGE-L (R-L), CIDEr to evaluate the performance of response generation. For human evaluation, we randomly sample 100 dialogues from testing set and employ crowdsourcing workers to rate generated responses based on five aspects of Empathy, Coherence, Informativity, Fluency, and Specificity. The score is from 1 to 5 (1: not at all, 3: OK, 5: very good), except Specificity. The Specificity score is 1 or 0, representing yes or no. Fleiss' Kappa of the human evaluation results is 0.498, indicating moderate agreement.

**Baselines & Hyper-parameters.** We choose MoEL [7], MIME [8], EmpDG [9], EC (soft) [11], KEMP [13], CEM [12], and DialoGPT (345M) [30] as baselines. For vertices in the graph, we use VGAE [31] to initialize representations, and the embedding size is 128. The hidden size of GRU is 768, and the maximum number of concepts is 5. We use Adam for optimization with the initial learning rate of 0.001.

### 3.2. Results and Analysis

**Automatic and Human Evaluations.** Table 1 reports the automatic and human experimental results. We observe that ECTG considerably exceeds baselines in most metrics for the automatic evaluation, demonstrating that ECTG is beneficial for empathetic dialogue generation. ECTG also achieves the best performance in four aspects for the human evaluation except Fluency, which verifies that ECTG can generate more empathetic, coherent, informative, and specific responses with the guidance of emotion causes and the transition of concepts. Additionally, we note that there is no significant difference in Fluency between models, and we speculate that the responses generated by all models are already fluent.

| Models | Automatic Evaluation | | | | | | Human Evaluation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Multi-trs | 2.103 | 0.1948 | 0.456 | 1.947 | 16.67 | 12.81 | 2.91 | 2.87 | 2.48 | 4.86 | 0.24 |
| MoEL | 1.933 | 0.2166 | 0.469 | 2.155 | 17.00 | 14.60 | 2.89 | 2.87 | 2.46 | 4.93 | 0.21 |
| MIME | 1.894 | 0.2039 | 0.449 | 1.829 | 16.64 | 13.68 | 3.13 | 2.97 | 2.59 | 4.89 | 0.24 |
| EmpDG | 1.975 | 0.2188 | 0.470 | 1.981 | 17.34 | 14.70 | 3.00 | 2.97 | 2.55 | 4.93 | 0.24 |
| EC (soft) | 1.345 | 0.1925 | 1.698 | 8.493 | 15.67 | 10.21 | 2.96 | 3.00 | 2.53 | 4.92 | 0.27 |
| KEMP | 1.762 | 0.1948 | 0.660 | 3.074 | 15.43 | 12.78 | 2.78 | 2.72 | 2.46 | **4.94** | 0.21 |
| CEM | 1.629 | 0.2134 | 0.645 | 2.856 | 16.27 | 15.83 | 3.02 | 3.21 | 2.38 | 4.90 | 0.19 |
| DialoGPT | 0.734 | 0.1515 | **3.140** | **17.551** | 8.51 | 7.00 | 3.70 | 3.89 | 3.06 | 4.88 | 0.61 |
| ECTG | **5.467** | **0.2701** | 1.840 | 16.404 | **23.77** | **51.43** | **3.78**‡ | **4.13**‡ | **3.13**‡ | 4.88 | **0.64**‡ |

**Table 1**. Automatic and human evaluations. †, ‡ represent the statistical significance (t-test) with p-value $<0.05$ and $0.01$.

| Models | B-4 | $F_{BERT}$ | Dist-1 | Dist-2 | R-L | CIDEr |
|---|---|---|---|---|---|---|
| ECTG | **5.47** | **0.2701** | 1.84 | **16.40** | **23.77** | **51.43** |
| w/o copy | 2.75 | 0.2569 | **2.37** | 14.73 | 20.95 | 39.62 |
| w/o seca | 3.04 | 0.2539 | 2.27 | 14.22 | 21.32 | 40.73 |
| w/o graph | 3.21 | 0.2301 | 2.18 | 13.56 | 18.88 | 23.98 |

**Table 2**. Results of the ablation study.

**Ablation Study.** We designed three variants of ECTG for the ablation study: **1) w/o copy**. We remove the Transformer decoder with the copy mechanism and only employ the non-autoregressive generation. **2) w/o seca**. The span-level emotion cause analysis is removed, then all keywords in the utterance are adopted to construct the graph. **3) w/o graph**. We remove the emotion cause transition graph and replace it with the form of text. The obtained results are shown in Table 2. We can observe that variants drop dramatically in most metrics, indicating our model settings' effectiveness. According to statistics, responses generated by ECTG tend to be longer than those generated by variants. It may have a great impact when calculating uni-gram. However, other metrics help prove that responses generated by ECTG are better.

**Case Study.** In Table 3, we provide some cases to compare generated responses of ECTG and baselines. In the first case, affective empathy oriented baselines roughly perceive the user's emotion status and respond generally. Although models with additional knowledge convey more information, their responses are not targeted to the context. EC(soft) successfully identifies the user's emotional state and replies with specific examples. However, the response is not particularly coherent due to the lack of global graph guidance. In contrast, ECTG understands the user's emotions and experiences accurately and gives good wishes with empathetic, relevant, and non-universal responses. In the second case of multiple-turn dialogue context, compared with other models that acknowledge the user's emotion, ECTG expresses appropriate emotion and explores more valuable information.

**Exploration Experiment.** We further explore the transferability of ECTG concepts by integrating predicted response concepts into the pre-trained model as prompts. We adopt a large-scale dialogue model DialoGPT (345M), whose parameter number is significantly higher than our model. We also choose BlenderBot [32] as the reference for the pre-trained model in the field of empathetic dialogue, which is trained with multiple communication skills. The results in Table 4 show

| Emotion | Nostalgic |
|---|---|
| Context | "I recently **spoke** with my **ex-girlfriend** on the phone. The **conversation** went pretty **well**, and it **reminded** me of my **past experiences** with her." |

| Multi-trs | That is so sweet. I am glad you were able to see her. |
|---|---|
| MoEL | That is so nice of you. I am glad you are okay. |
| MIME | That is nice of you. I am sure you miss her a lot of hard work. |
| EmpDG | That is very nice of you. |
| EC (soft) | Oh, that's so sweet. I'm sure you were very grateful for her. |
| KEMP | That is great. I am sure you will have a great time to have a blast. |
| CEM | That is great. Did you find out of it? |
| DialoGPT | I hope everything turned out well for you. |

| ECTG | That's <span style="color:red">interesting</span>. I <span style="color:red">hope</span> it turns out well for you two can still be <span style="color:red">together</span>. |
|---|---|

| Emotion | Surprised |
|---|---|
| Context | "I couldnt celebrate my 18th birthday as I had exams during the time and I felt terrible for it. But my **friends** threw a **surprise party** for me and it was the best!" "That's **great**! How did they **surprise** you?" "I was studying in my bedroom for a couple of hours straight with headphones on and then my mom called me to the living room for dinner and when I went out the **whole place** was **decorated** and **everyone was there**." |

| Multi-trs | That is good. I am glad you are able to do that. |
|---|---|
| MoEL | That is so nice of them! |
| MIME | That is a good thing to do. I am sure you are going to have a good time. |
| EmpDG | That is good. I am sure you will be able to get a little one. |
| EC (soft) | That's awesome! I hope you get the next step! |
| KEMP | I am sure you will do great! |
| CEM | That is great! I love the day! |
| DialoGPT | That's so nice of them! |

| ECTG | That's so <span style="color:red">sweet</span> of them! Did you <span style="color:red">enjoy</span> it? |
|---|---|

**Table 3**. Cases of responses generated by all models. The identified spans are marked with shadow. The concepts in the context are bold, and the predicted concepts are colored red.

that DialoGPT with concepts of ECTG outperforms DialoGPT and BlenderBot in most metrics, which verifies that combining predicted response concepts can improve performance.

| Models | B-4 | $F_{BERT}$ | Dist-1 | Dist-2 | R-L | CIDEr |
|---|---|---|---|---|---|---|
| BlenderBot | 1.30 | **0.1599** | 2.41 | 12.45 | 9.71 | 8.07 |
| DialoGPT | 0.73 | 0.1515 | 3.14 | 17.55 | 8.51 | 7.00 |
| + ECTG concepts | **1.36** | 0.1524 | **3.75** | **21.39** | **10.78** | **15.91** |

**Table 4**. Results of the exploration experiment.

## 4. CONCLUSION

In this paper, we propose to generate empathetic responses aware of emotion cause concepts. We construct an emotion cause transition graph to explicitly model natural transitions in the human empathetic dialogue and design a model using the graph to benefit the empathetic response generation. Automatic and human evaluations verify our approach's ability in the field of empathetic dialogue.

# 5. REFERENCES

[1] Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum, "The design and implementation of xiaoice, an empathetic social chatbot," *Comput. Linguistics*, vol. 46, no. 1, pp. 53–93, 2020.

[2] Jan Wieseke, Anja Geigenmüller, and Florian Kraus, "On the role of empathy in customer-employee interactions," *Journal of service research*, vol. 15, no. 3, pp. 316–331, 2012.

[3] Ting-En Lin, Yuchuan Wu, Fei Huang, Luo Si, Jian Sun, and Yongbin Li, "Duplex conversation: Towards human-like interaction in spoken dialogue systems," in *KDD*, 2022.

[4] Wanwei He, Yinpei Dai, Yinhe Zheng, Yuchuan Wu, Zheng Cao, Dermot Liu, Peng Jiang, Min Yang, Fei Huang, Luo Si, et al., "Space: A generative pre-trained model for task-oriented dialog with semi-supervised learning and explicit policy injection," *AAAI*, 2022.

[5] Verónica Pérez-Rosas, Rada Mihalcea, Kenneth Resnicow, Satinder Singh, and Lawrence C. An, "Understanding and predicting empathic behavior in counseling therapy," in *ACL*, 2017.

[6] Meghan L Healey and Murray Grossman, "Cognitive and affective perspective-taking: evidence for shared and dissociable anatomical substrates," *Frontiers in neurology*, vol. 9, 2018.

[7] Zhaojiang Lin, Andrea Madotto, Jamin Shin, Peng Xu, and Pascale Fung, "Moel: Mixture of empathetic listeners," in *EMNLP-IJCNLP*, 2019, pp. 121–132.

[8] Navonil Majumder, Pengfei Hong, Shanshan Peng, Jiankun Lu, Deepanway Ghosal, Alexander F. Gelbukh, Rada Mihalcea, and Soujanya Poria, "MIME: mimicking emotions for empathetic response generation," in *EMNLP*, 2020, pp. 8968–8979.

[9] Qintong Li, Hongshen Chen, Zhaochun Ren, Pengjie Ren, Zhaopeng Tu, and Zhumin Chen, "Empdg: Multi-resolution interactive empathetic dialogue generation," in *COLING*, 2020.

[10] Hyunwoo Kim, Byeongchang Kim, and Gunhee Kim, "Perspective-taking and pragmatics for generating empathetic responses focused on emotion causes," in *EMNLP*, 2021.

[11] Jun Gao, Yuhan Liu, Haolin Deng, Wei Wang, Yu Cao, Jiachen Du, and Ruifeng Xu, "Improving empathetic response generation by recognizing emotion cause in conversations," in *EMNLP*, 2021, pp. 807–819.

[12] Sahand Sabour, Chujie Zheng, and Minlie Huang, "Cem: Commonsense-aware empathetic response generation," in *AAAI*, 2022, vol. 36, pp. 11229–11237.

[13] Qintong Li, Piji Li, Zhaochun Ren, Pengjie Ren, and Zhumin Chen, "Knowledge bridging for empathetic dialogue generation," in *AAAI*, 2022, pp. 10993–11001.

[14] Anna Wierzbicka, *Emotions across languages and cultures: Diversity and universals*, Cambridge university press, 1999.

[15] Guimin Hu, Ting-En Lin, Yi Zhao, Guangming Lu, Yuchuan Wu, and Yongbin Li, "Unimse: Towards unified multimodal sentiment analysis and emotion recognition," *arXiv preprint arXiv:2211.11256*, 2022.

[16] Sai Zhang, Yuwei Hu, Yuchuan Wu, Jiaman Wu, Yongbin Li, Jian Sun, Caixia Yuan, and Xiaojie Wang, "A slot is not built in one utterance: Spoken language dialogs with sub-slots," in *Findings of ACL*, 2022, pp. 309–321.

[17] Yinpei Dai, Wanwei He, Bowen Li, Yuchuan Wu, Zheng Cao, Zhongqi An, Jian Sun, and Yongbin Li, "Cgodial: A large-scale benchmark for chinese goal-oriented dialog evaluation," *arXiv preprint arXiv:2211.11617*, 2022.

[18] Yinpei Dai, Hangyu Li, Yongbin Li, Jian Sun, Fei Huang, Luo Si, and Xiaodan Zhu, "Preview, attend and review: Schema-aware curriculum learning for multi-domain dialogue state tracking," in *EMNLP*, 2021, pp. 879–885.

[19] Houyu Zhang, Zhenghao Liu, Chenyan Xiong, and Zhiyuan Liu, "Grounded conversation generation as guided traverses in commonsense knowledge graphs," in *ACL*, 2020.

[20] Yicheng Zou, Zhihua Liu, Xingwu Hu, and Qi Zhang, "Thinking clearly, talking fast: Concept-guided non-autoregressive generation for open-domain dialogue systems," in *EMNLP*, 2021.

[21] Xiangju Li, Wei Gao, Shi Feng, Yifei Zhang, and Daling Wang, "Boundary detection with BERT for span-level emotion cause analysis," in *Findings of ACL*, 2021, pp. 676–682.

[22] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy, "Spanbert: Improving pre-training by representing and predicting spans," *Trans. Assoc. Comput. Linguistics*, vol. 8, pp. 64–77, 2020.

[23] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly, "Pointer networks," in *Advances in Neural Information Processing Systems*, 2015, vol. 28.

[24] Jianheng Tang, Tiancheng Zhao, Chenyan Xiong, Xiaodan Liang, Eric P. Xing, and Zhiting Hu, "Target-guided open-domain conversation," in *ACL*, 2019, pp. 5624–5634.

[25] Wanwei He, Yinpei Dai, Binyuan Hui, Min Yang, Zheng Cao, Jianbo Dong, Fei Huang, Luo Si, and Yongbin Li, "Space-2: Tree-structured semi-supervised contrastive pre-training for task-oriented dialog understanding," *arXiv preprint arXiv:2209.06638*, 2022.

[26] Wanwei He, Yinpei Dai, Min Yang, Jian Sun, Fei Huang, Luo Si, and Yongbin Li, "Space-3: Unified dialog model pre-training for task-oriented dialog understanding and generation," *arXiv preprint arXiv:2209.06664*, 2022.

[27] Mitchell Stern, William Chan, Jamie Kiros, and Jakob Uszkoreit, "Insertion transformer: Flexible sequence generation via insertion operations," in *ICML*, 2019, pp. 5976–5985.

[28] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau, "Towards empathetic open-domain conversation models: A new benchmark and dataset," in *ACL*, 2019, pp. 5370–5381.

[29] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi, "Bertscore: Evaluating text generation with BERT," in *ICLR*, 2020.

[30] Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan, "DIALOGPT : Large-scale generative pre-training for conversational response generation," in *ACL*, 2020.

[31] Thomas N. Kipf and Max Welling, "Variational graph auto-encoders," *CoRR*, vol. abs/1611.07308, 2016.

[32] Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston, "Recipes for building an open-domain chatbot," in *EACL*, 2021, pp. 300–325.