

# TREA: Tree-structure Reasoning Schema for Conversational Recommendation

Wendi Li<sup>1,2</sup>, Wei Wei<sup>1,2,✉</sup>, Xiaoye Qu<sup>1</sup>,

Xianling Mao<sup>3</sup>, Ye Yuan<sup>4</sup>, Wenfeng Xie<sup>4</sup>, Dangyang Chen<sup>4</sup>

<sup>1</sup>Cognitive Computing and Intelligent Information Processing (CCIIP) Laboratory,  
Huazhong University of Science and Technology

<sup>2</sup>Joint Laboratory of HUST and Pingan Property & Casualty Research (HPL)

<sup>3</sup>Department of Computer Science and Technology, Beijing Institute of Technology

<sup>4</sup>Ping An Property & Casualty Insurance company of China

<sup>1</sup>{wendili, weiw, xiaoye}@hust.edu.cn <sup>3</sup>maoxl@bit.edu.cn

<sup>4</sup>{yuanye503, xiewenfeng801, chendangyang273}@pingan.com.cn

## Abstract

Conversational recommender systems (CRS) aim to timely trace the dynamic interests of users through dialogues and generate relevant responses for item recommendations. Recently, various external knowledge bases (especially knowledge graphs) are incorporated into CRS to enhance the understanding of conversation contexts. However, recent reasoning-based models heavily rely on simplified structures such as linear structures or fixed-hierarchical structures for causality reasoning, hence they cannot fully figure out sophisticated relationships among utterances with external knowledge. To address this, we propose a novel Tree-structure Reasoning schEmA named TREA. TREA constructs a multi-hierarchical scalable tree as the reasoning structure to clarify the causal relationships between mentioned entities, and fully utilizes historical conversations to generate more reasonable and suitable responses for recommended results. Extensive experiments on two public CRS datasets have demonstrated the effectiveness of our approach. Our code is available at <https://github.com/WindyLee0822/TREA>

## 1 Introduction

Conversation Recommender System (CRS) has become increasingly popular as its superiority in timely discovering user dynamic preferences in practice. As opposed to traditional passive-mode recommendation systems, it highlights the importance of proactively clarifying and tracing user interests through live conversation interactions, which notably enhance the success rate of item recommendations.

Since sole contextual utterances are insufficient for comprehensively understanding user preferences, there are many efforts devoted to incorporat-

ing various external knowledge (Chen et al., 2019; Zhou et al., 2020a, 2022; Wang et al., 2022; Yang et al., 2022), which typically enrich the contextual information with mentioned entities recognized over utterances. However, these methods fail to model the complex causal relations among mentioned entities, owing to the diversity of user interest expression and the frequent shift of conversation topic as shown in Figure 1.

Actually, it is non-trivial to explicitly model the complex causal relationships of conversations. Although there are several reasoning-based methods proposed for CRS, their simplified structures make the objective unattainable. Some researches (Zhou et al., 2021) track the mentioned entities as linear sequential fragments analogous to (1) in Figure 1. However, the linear structure is only suitable for adjacent relation modeling, which may not always work well since the actual causality between mentioned entities exists multi-hop jumps ("comedy"- "La La Land" in Figure 1). Other studies (Ma et al., 2021) propose other forms of specially-designed structures for reasoning akin to (2) in Figure 1, but they generally have fixed hierarchies, which often degenerate into a simple 2-layer hierarchy "history"- "prediction", neglecting the causal relations of historical entities. Therefore, neither of them is applicable for full modeling of the complex reasoning causality within conversations.

To improve the reasoning capability of CRS, the challenges are twofold. The first challenge lies in empowering the model to illuminate the causal inference between all mentioned entities. To tackle this, we perform abductive reasoning for each mentioned entity to construct the multi-hierarchical reasoning tree. The reasoning tree explicitly preserves logical relations between all entities and can be continuously expanded as the conversation continues, which provides the model with a clear

✉ Corresponding Author

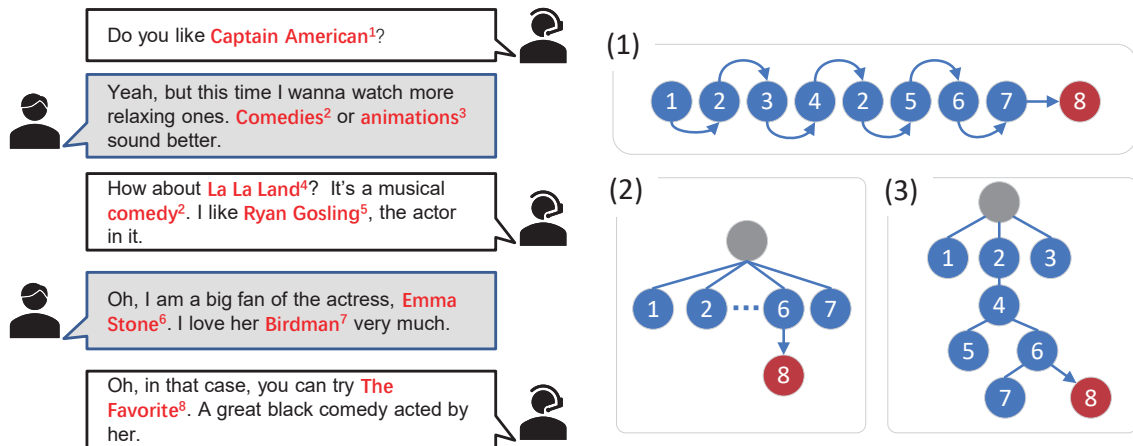


Figure 1: An example of conversational recommendation scenarios and three kinds of reasoning structures for CRS. In the conversation example, entities are marked in red and the upper-left number corresponds to the figure in the reasoning structure. (1) corresponds to the linear structure. (2) corresponds to the structure with two fixed hierarchies (history-prediction), flattening all the mentioned entities at the first hierarchy. (3) corresponds to our multi-hierarchical structure of TREA.

reference to historical information for prediction. The second challenge is how to utilize reasoning information in response generation. We enable the model to extract relevant textual information from the historical conversation with the corresponding reasoning branch, thus promoting the correlation between generated responses and recommended items. We name this **Tree-structure Reasoning schEmA TREA**.

To validate the effectiveness of our approach, we conduct experiments on two public CRS datasets. Experimental results show that our TREA outperforms competitive baselines on both the recommendation and conversation tasks. Our main contributions are summarized as follows:

- To the best of our knowledge, it is the first trial of CRS to reason every mentioned entity for its causation.
- We propose a novel tree-structured reasoning schema to clarify the causality relationships between entities and mutual the reasoning information with the generation module.
- Extensive experiments demonstrate the effectiveness of our approach in both the recommendation and conversation tasks.

## 2 Related Work

Conversational Recommender System (CRS) explores user preference through natural language

dialogues. Previous works can be roughly categorized into two types. The first category of CRS is recommendation-biased CRS (Sun and Zhang, 2018; Lei et al., 2020b,a; Deng et al., 2021; Zhang et al., 2022). This category focuses solely on interactive recommendations but the function of natural language is ignored. Several fixed response templates are preset on the agents and users cannot use free text but only have limited options, which can be detrimental to the user experience.

The other category of CRSs is dialog-biased CRS (Li et al., 2018; Moon et al., 2019; Chen et al., 2020; Liu et al., 2021; Sarkar et al., 2020). This category emphasizes the critical effect of natural language, aiming to understand user utterances for accurate recommendations and generate human-like responses. Noticing that entities (Gu et al., 2022; Qu et al., 2022, 2023) mentioned in conversations are important cues for modeling user preferences, Chen et al. (2019) firstly integrates KG to enhance the user representation. Zhou et al. (2020a); Liang et al. (2021) use two KGs on entity-granularity and word-granularity respectively to represent the user preference more comprehensively. Subsequent researches introduce other types of external knowledge e.g. item description (Lu et al., 2021; Zhou et al., 2022) or pretrained language models (PLMs) (Yang et al., 2022; Wang et al., 2022) to further assist the user representations. However, they commonly treat each mentioned knowledge piece equally and integrate them into an aggregated representation.

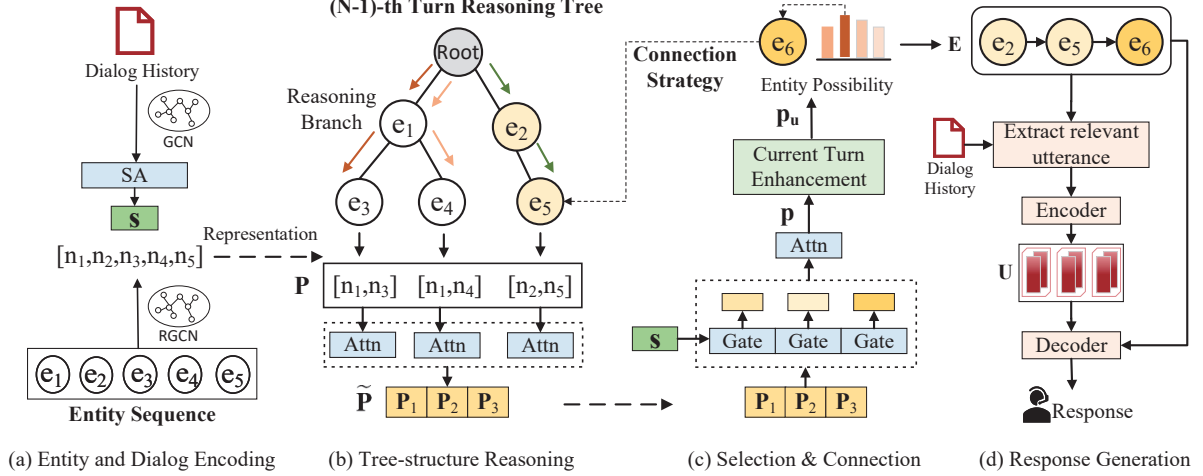


Figure 2: The overview of our proposed TREA. We first encode the entities and the sentences in the dialog history. Then we aggregate the information of each reasoning branch in the current reasoning tree. Later, a comprehensive representation of dialog semantics measures the devotion of each reasoning branch to the current recommendation. After the current turn enhancement, we select the entity to join the reasoning tree with the connection strategy. The extended reasoning branch guide the extraction of relevant textual information for the generation module.

Recently, some researches manage to model the reasoning process during conversations. Zhou et al. (2021) linearize the mentioned entity sequence and reasoning the inferential causality between the adjacent entity pairs. Ma et al. (2021) create non-linear reasoning structures, but they do not preserve the hierarchy of historical turns. Therefore these reasoning methods have limited performance improvement.

To sort out the causal relations among utterances, our model performs tree-structured reasoning on the entire dialogue history for each mentioned entity. We also inject the reasoning information into the generation process to make responses more relevant, achieving that the reasoning process facilitates both recommendation and generation tasks simultaneously.

### 3 Methods

In this section, we present the Tree-structure reasoning schema TREA as demonstrated in Figure 2. Specifically, we first introduce the encoding of entities and word tokens. Then we illustrate the construction procedure of the reasoning tree. Later, we describe how the reasoning information supports the generation module. Finally, we explain the process of parameter optimization.

#### 3.1 Entity and Dialog Encoding

Following previous works (Chen et al., 2019; Zhou et al., 2020a; Ma et al., 2021; Zhou et al., 2022),

we first perform entity linking based on an external KG DBpedia (Bizer et al., 2009), and then encode the relational semantics via a relational graph neural network (RGCN) (Schlichtkrull et al., 2018) to obtain the corresponding entity embeddings. Formally, the embedding  $\mathbf{n}_e^{l+1}$  of entity  $e$  at the  $l+1$ -th graph layer is calculated as:

$$\mathbf{n}_e^{l+1} = \sigma\left(\sum_{r \in R} \sum_{e' \in \mathcal{N}_e^r} \frac{1}{Z_{e,r}} \mathbf{W}_r^l \mathbf{n}_{e'}^l + \mathbf{W}^l \mathbf{n}_e^l\right) \quad (1)$$

where  $R$  is a relation set,  $\mathcal{N}_e^r$  denotes the set of neighboring nodes for  $e$  under the relation  $r$ ,  $\mathbf{W}_r^l$ ,  $\mathbf{W}^l$  are learnable matrices for relation-specific aggregation with neighboring nodes and representation transformation respectively,  $Z_{e,r}$  is a normalization factor,  $\sigma$  denotes the sigmoid function. The semantic information of word tokens is encoded by an external lexical knowledge graph ConceptNet (Speer et al., 2017). We further adopt a graph convolutional neural network (GCN) (Kipf and Welling, 2016) to propagate and aggregate information over the entire graph.

#### 3.2 Reasoning Tree Construction.

The construction of reasoning trees is introduced in a manner similar to mathematical induction. We first explain the structure initialization at the first conversation round, then illustrate the structure transition from the  $(n-1)$ -th round to the  $n$ -th round. The structure of the whole tree can be deduced accordingly.

To initialize the reasoning tree, we first set a pseudo node as the root node. The root node does not represent any entity in the conversations but is just a placeholder. When the first utterance is coming, the first mentioned entity is directly connected to the root node. The subsequent entities in the first utterance are connected following the Algorithm 1.

When the conversation progresses to  $(n-1)$ -th round, the known conditions are as follows: the current reasoning tree  $\mathcal{T}_{n-1}$ , utterance tokens sequences  $s_t$ . They are utilized for the extension of the reasoning tree  $\mathcal{T}_{n-1}$ , which is described in two parts, tree-structure reasoning and the selection & connection of candidate entities.

**Tree-Structure Reasoning.** We embed all the reasoning branches and pad them to a certain length  $l_r$ . A path from the root node to any leaf node of the tree is referred to as a *reasoning branch* since it expresses a chain of coherent inferences. To represent the sequential information for each reasoning branch, we inject a learnable position embedding into the embedding of each entity element. The position-enhanced branch embedding matrix is denoted as  $\mathbf{P} \in \mathbb{R}^{n_r \times l_r \times d}$  where  $n_r$  is the branch number of  $\mathcal{T}_{n-1}$  and  $d$  is the dimension of embeddings. We incorporate a linear attention mechanism to integrate the representation of each path. The attention scores are calculated as follows:

$$\begin{aligned} \tilde{\mathbf{P}} &= \text{Attn}(\mathbf{P}) = \mathbf{P} \alpha_r \\ \alpha_r &= \text{Softmax}(\mathbf{b}_r \tanh(\mathbf{W}_r \mathbf{P})) \end{aligned} \quad (2)$$

where  $\mathbf{W}_r, \mathbf{b}_r$  are learnable parameters. Embeddings of entities in a certain reasoning branch are aggregated according to the attention score. Then we can obtain the comprehensive representations of reasoning branches denoted as  $\tilde{\mathbf{P}} \in \mathbb{R}^{n_r \times d}$ .

**Selection & Connection.** Since the reasoning branches have varying-degrees contributions to the next-hop entity, the model analyzes the semantics of word tokens  $s_t$  to measure the impact of each branch. The formulas are as follows:

$$\begin{aligned} \mathbf{p} &= \text{Attn}(\gamma \tilde{\mathbf{P}} + (1 - \gamma) \mathbf{s}) \\ \gamma &= \sigma(\mathbf{W}_s \text{Concat}(\tilde{\mathbf{P}}, \mathbf{s})) \end{aligned} \quad (3)$$

where  $\mathbf{W}_s$  is a learnable parameter,  $\mathbf{s}$  is the comprehensive semantic representation of the word tokens in ConceptNet which are aggregated with the linear attention mechanism in Eq.2. Then we can obtain the user representation  $\mathbf{p}_u$  that combines semantic and reasoning information. Since the latest turn has a prominent significance to the response (Li et al.,

2022), we collect the entities and word tokens from the current conversation turn, embedded to  $\mathbf{e}_c, \mathbf{s}_c$ . Then we aggregate the current turn information and mutual it with acquired representation  $\mathbf{p}$  as follows:

$$\mathbf{p}_u = g(\mathbf{p}, g'(\text{Attn}(\mathbf{e}_c), \text{Attn}(\mathbf{s}_c))) \quad (4)$$

where  $g(\cdot, \cdot), g'(\cdot, \cdot)$  are two gate layers like Eq.3. Then we derive the next-hop possibility distribution from the overall user representation:

$$\mathcal{P}_r^u = \text{Softmax}([\mathbf{p}_u \mathbf{e}_0^T, \dots, \mathbf{p}_u \mathbf{e}_n^T]) \quad (5)$$

where  $\mathbf{e}_0, \dots, \mathbf{e}_n$  are representations of all entities. The entity with the largest probability is selected and connected to the reasoning tree. The connection strategy is as Algorithm 1.

---

#### Algorithm 1: Connection Strategy

---

**input:** Selected entity  $e^*$ ; Entity sequence  $ES$  in reverse order of mention;  
Reasoning Tree  $\mathcal{T}$  with root node  $r$

```

1 foreach  $e$  in  $ES$  do
2   if  $\text{IsAdj}(e, e^*)$  then
3     // Two entities are adjacent in KG;
4      $\text{AddEdge}(e, e^*)$ ;
5     // Add an edge  $(e, e^*)$  in  $\mathcal{T}$ ;
6     return
7   end
8 end
9  $\text{AddEdge}(r, e^*)$ ;
10 return

```

---

### 3.3 Reasoning-guided Response Generation

After adding the predicted entity to the reasoning tree, the objective of the conversation module is to generate utterances with high relevance to the predicted entity. Reasoning branches that involve the new entity and the historical utterances that mention the relevant entities in branches are extracted, which are encoded by RGCN and standard Transformer (Vaswani et al., 2017) respectively. The corresponding embedding matrices are denoted as  $\mathbf{E}, \mathbf{U}$ . Following (Zhou et al., 2020a), we incorporate multiple cross-attention layers in a Transformer-variant decoder to fuse the two groups of information. The probability distribution over the vocabulary is calculated as follows:

$$\mathbf{R}^l = \text{Decoder}(\mathbf{R}^{l-1}, \mathbf{E}, \mathbf{U}) \quad (6)$$

$$\mathbf{R}^b = \text{FFN}(\text{Concat}(\text{Attn}(\mathbf{E}), \mathbf{R}^l)) \quad (7)$$

$$\mathcal{P}_g = \text{Softmax}(\mathbf{R}^l \mathbf{V}^T + \mathbf{R}^b \mathbf{W}^v) \quad (8)$$



where  $\mathbf{V}$  is the embedding matrix of all words in the vocabulary,  $\mathbf{W}^v$  is a learnable parameter that converts the  $\mathbf{R}^b$  dimension to  $|\mathbf{V}|$ . The copy mechanism is adopted in Eq.7 to enhance the generation of knowledge-related words. The transformation chain (Zhou et al., 2020a) in the decoder of Eq.6 is *generated words*  $\rightarrow$  *relevant entities*  $\rightarrow$  *historical utterances*.

### 3.4 Optimization

The parameters can be categorized into two parts, the reasoning parameters and the generation parameters, denoted by  $\theta_r, \theta_g$ . The reasoning objective is to maximize the predicted probability of the upcoming entity. The cross-entropy loss is adopted to train the reasoning module. During the training, we propose two auxiliary loss functions, isolation loss to maintain the independence of each reasoning branch, and alignment loss to bridge the representation gap.

**Isolation Loss.** Since reasoning branches that have no shared parts are generally irrelevant, representations from different reasoning branches are expected to be dissimilar. To maintain the isolation of each reasoning branch, we propose isolation loss. Given representations of different reasoning branches, the isolation loss is calculated as

$$\mathcal{L}_I = \sum_{i \neq j} \text{sim}(\tilde{\mathbf{p}}_i, \tilde{\mathbf{p}}_j) = \sum_{i \neq j} \frac{\tilde{\mathbf{p}}_i \tilde{\mathbf{p}}_j}{|\tilde{\mathbf{p}}_i| \cdot |\tilde{\mathbf{p}}_j|} \quad (9)$$

where  $\tilde{\mathbf{p}}_i, \tilde{\mathbf{p}}_j$  are representations of two different reasoning branches extracted from  $\tilde{\mathbf{P}}$ .

**Alignment Loss.** The representation gap exists between the semantics and the entities since their encoding processes are based on two separate networks. Hence the entity representation and semantic representation of the same user should be dragged closer; those of different users should be pushed further to reduce the gap. The formula is as follows:

$$\mathcal{L}_a = \lambda_c \text{sim}(\mathbf{p}_c, \mathbf{s}_c) + (1 - \lambda_c) \text{sim}(\mathbf{p}, \mathbf{s}) \quad (10)$$

where  $\mathbf{p}_c, \mathbf{s}_c$  are aggregated representation  $\text{Attn}(\mathbf{e}_c), \text{Attn}(\mathbf{w}_c)$  in Eq.4,  $\lambda_c$  is a hyperparameter.

Then We can optimize parameters  $\theta_r$  through the following formula:

$$\mathcal{L}_r = - \sum_u \sum_{e_i} \log \mathcal{P}_r^u[e_i] + \lambda_I \mathcal{L}_I + \lambda_a \mathcal{L}_a \quad (11)$$

where  $e_i$  is the order of the target entity at the  $i$ -th conversation round of user  $u$ ,  $\lambda_I, \lambda_a$  are hyperparameters.

When the reasoning loss  $\mathcal{L}_r$  converges, we optimize the parameters in  $\theta_g$ . After obtaining the relevant entities and utterances via the reasoning tree, we calculate the probability distribution of the next token. To learn the generation module, we set the cross-entropy loss as:

$$\mathcal{L}_g = - \frac{1}{N} \sum_{t=1}^N \log \mathcal{P}_g^t(s_t | s_1, s_2, \dots, s_{t-1}) \quad (12)$$

where  $N$  is the number of turns in a certain conversation  $C$ . We compute this loss for each utterance  $s_t$  from  $C$ .

## 4 Experiment

### 4.1 Dataset.

We conduct our experiments on two widely-applied benchmark datasets on CRS, which are multilingual including English (ReDial) and Chinese (TG-ReDial). **ReDial**(Li et al., 2018) collects high-quality dialogues for recommendations on movies through crowd-sourcing workers on Amazon Mechanical Turk(AMT). The workers create conversations for the task of movie recommendation in a user-recommender pair setting following a set of detailed instructions. It contains 10,006 conversations consisting of 182,150 utterances. **TG-ReDial**(Zhou et al., 2020b) is annotated in a semi-automatic way. It emphasizes natural topic transitions from non-recommendation scenarios to the desired recommendation scenario. Each conversation includes a topic path to enforce natural semantic transitions. It contains 10,000 conversations consisting of 129,392 utterances.

### 4.2 Baselines

We evaluate the effectiveness of our model with following competitive baselines:

*ReDial* (Li et al., 2018) comprises a conversation module based on hierarchical encoder-decoder architecture(Serban et al., 2017) and a recommendation module based on auto-encoder.

*KBRD* (Chen et al., 2019) firstly utilizes KG to enhance the user representation. The Transformer(Vaswani et al., 2017) architecture is applied in the conversation module.

*KGSF* (Zhou et al., 2020a) incorporate two external knowledge graphs on different aspects to further enhance the user representations. The KG information is employed in the decoding process.

Dataset	ReDial						TG-ReDial					
Method	R@10	R@50	Dist-3	Dist-4	Bleu-2	Bleu-3	R@10	R@50	Dist-3	Dist-4	Bleu-2	Bleu-3
ReDial	0.140	0.320	0.269	0.464	0.022	0.008	0.002	0.013	0.529	0.801	0.041	0.010
KBRD	0.150	0.336	0.288	0.489	0.024	0.009	0.032	0.077	0.691	0.997	0.042	0.012
KGSF	0.183	0.377	0.302	0.518	0.025	0.009	0.030	0.074	1.045	1.579	0.046	0.014
RevCore	0.204	0.392	0.307	0.528	0.025	0.010	0.029	0.075	1.093	1.663	0.047	0.014
CR-Walker	0.187	0.373	0.338	0.557	0.024	0.009	-	-	-	-	-	-
CRFR	0.202	0.399	0.516	0.639	-	-	-	-	-	-	-	-
C <sup>2</sup> -CRS	0.208	0.409	0.412	0.622	0.027	0.012	0.032	0.078	1.210	1.691	0.048	0.015
UCCR	0.202	0.408	0.329	0.564	0.026	0.011	0.032	0.075	1.197	1.668	0.049	0.014
<b>TREA</b>	<b>0.213*</b>	<b>0.416*</b>	<b>0.692*</b>	<b>0.839*</b>	<b>0.028*</b>	<b>0.013*</b>	<b>0.037*</b>	<b>0.110*</b>	<b>1.233*</b>	<b>1.712*</b>	<b>0.050*</b>	<b>0.017*</b>

Table 1: Automatic evaluation results on two datasets. Boldface indicates the best results. Significant improvements over best baseline marked with \*. (t-test with  $p < 0.05$ )

*CRFR* (Zhou et al., 2021) can generate several linear reasoning fragments through reinforcement learning to track the user preference shift.

*CR-Walker* (Ma et al., 2021) create a two-hierarchy reasoning tree between history and forecast and preset several dialog intents to guide the reasoning.

*C<sup>2</sup>-CRS* (Zhou et al., 2022) proposed a contrastive learning based pretraining approach to bridge the semantic gap between three external knowledge bases.

*UCCR* (Li et al., 2022) considers multi-aspect information from the current session, historical sessions, and look-alike users for comprehensive user modeling.

### 4.3 Metrics

For recommendation evaluation, we used *Recall@n* ( $R@n, n=10,50$ ), which shows whether the top- $n$  recommended items include the ground truth suggested by human recommenders. For the response generation task, we evaluate models by *Bleu-n* ( $n=2,3$ ) (Papineni et al., 2002), *Dist-n* ( $n=3,4$ ) (Li et al., 2016) for word-level matches and diversity. To evaluate the generation performance more equitably, three annotators are invited to score the generated candidates from two datasets for human evaluation on the following three aspects: *Fluency*, *Relevance*, and *Informativeness*. The inter-annotator coherence is measured by Fleiss’ Kappa.

### 4.4 Implementation Details

We keep the same data preprocessing steps and hyperparameter settings as previous researches (Zhou et al., 2022; Ma et al., 2021). We adopt the same mask mechanism as NTRD (Liang et al., 2021).

The embedding dimensions of reasoning and generation are set to 300 and 128 respectively. In the encoding module, the word embeddings are initialized via Word2Vec<sup>1</sup> and the layer number is set to 1 for both GNN networks. The normalization constant of RGCN is 1. We use Adam optimizer (Kingma and Ba, 2015) with the default parameter setting. For training, the batch size is set to 64, the learning rate is 0.001, gradient clipping restricts the gradients within  $[0,0.02]$ . For hyperparameters,  $Z_e, r$  of RGCN in Eq.1 is 1,  $\lambda_c$  of representation alignment in Eq.10 is 0.9,  $\lambda_I, \lambda_a$  in Eq.11 is 0.008, 0.002 respectively.

### 4.5 Overall Performance Analysis

**Recommendation.** The columns R@10, R@50 of Table 1 present the evaluation results on the recommendation task. It shows that our TREA significantly outperforms all the baselines by a large margin on both two datasets, which verifies that TREA can clarify the sophisticated causality between the historical entities and accurately model the user preferences. Moreover, even though RevCore and C<sup>2</sup>-CRS utilize the additional knowledge, they are still not as effective as TREA, which further proves the significance of correct reasoning. CR-walker and CRFR are two previous methods that manage to reason over the background knowledge. CR-Walker does not preserve the hierarchy between the historical information and CRFR linearizes the reasoning structure. Therefore even though CR-walker conducts the additional annotations of dialog intents and CRFR applies the reasoning on another

<sup>1</sup><https://radimrehurek.com/gensim/models/word2vec.html>

KG to assist, the performance raising is limited, which certifies that our non-linear tree-structured reasoning over all mentioned entities does facilitate the user modeling.

Method	Rel.	Inf.	Flu.	Kappa
RevCore	1.98	2.22	1.53	0.78
CR-Walker	1.79	2.15	1.68	0.77
C <sup>2</sup> -CRS	2.02	2.25	1.69	0.66
UCCR	2.01	2.19	1.72	0.72
<b>TREA</b>	<b>2.43</b>	<b>2.26</b>	<b>1.75</b>	<b>0.75</b>

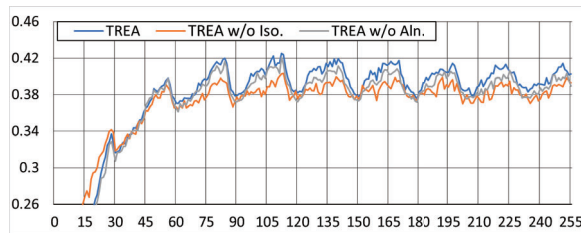
Table 2: Human evaluation results on the conversation task. Rel., Inf. and Flu. stand for Relevance, Informativeness and Fluency respectively. Boldface indicates the best results (t-test with  $p < 0.05$ ).

**Generation.** The columns Dist-n, Bleu-n of Table 1 present the automatic evaluation results on the conversation task. Since CR-walker adopts GPT-2 in the original model, we initialize the generation module with Word2Vec instead for a fair comparison. It shows that TREA surpasses all baselines on generation diversity and matchness. Table 2 presents the human evaluation results. All Fleiss’s kappa values exceed 0.6, indicating crowd-sourcing annotators have reached an agreement. The results show that our TREA leads to a higher relevance of generated utterances. It can be derived that the extraction of relevant information with the reasoning tree does improve the relevance of the generation.

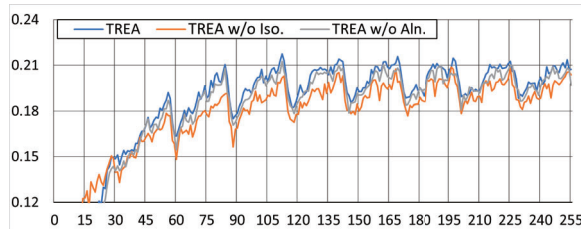
#### 4.6 Ablation Study

**Recommendation.** The parameter optimization for the reasoning module involves two additional loss, isolation loss (Iso.)  $\mathcal{L}_{\mathcal{I}}$  and alignment loss (Aln.)  $\mathcal{L}_a$ . We would like to verify the effectiveness of each part. We incorporate three variants of our model for ablation analysis on the recommendation task, namely *TREA w/o Iso.*, *TREA w/o Aln.* and *TREA w/o IA.*, which remove the isolation loss, the alignment loss and both of them respectively. As shown in Table 3, both components contribute to the final performance. Furthermore, we can see that removing the isolation loss leads to a large performance decrease, which suggests that maintaining the representation dependence of each reasoning branch is crucial to the correctness of the reasoning.

To further confirm that the performance improvement is consistent and stable instead of acciden-



(a) Recall@50



(b) Recall@10

Figure 3: Performance comparison of TREA and its two variants. One step (X-axis) denotes parameter updates for 20 batches of training data.

Dataset	ReDial		TG-ReDial	
Method	R@10	R@50	R@10	R@50
TREA	0.214	0.418	0.037	0.110
TREA w/o Iso.	0.202	0.405	0.028	0.079
TREA w/o Aln.	0.209	0.412	0.035	0.103
TREA w/o IA.	0.201	0.403	0.026	0.076

Table 3: Ablation results on the recommendation task. (t-test with  $p < 0.05$ )

tal. We test the models under different iteration steps and display the corresponding results in Figure 3. It can be seen that when the training loss converges, each ablation component contributes to the model performance regardless of the iteration number, which proves that the two additional loss functions are stably effective.

**The Effect of Isolation Loss.** The above subsection has verified the great impact of the isolation loss. We take a deeper dive to determine how it benefits model performance. If removing the isolation loss, since each reasoning branch participates in the calculation of the predicted possibility distribution, the representations of entities in different reasoning branches would approach each other for sharper descending of the loss value, which means that the representation of irrelevant entities would get similar irrationally and finally lead to the representation convergence of the entire knowledge graph. To confirm the assumption, we display the entity embeddings trained by TREA and TREA



Figure 4: 2D projection of KG embeddings trained by TREA (the above) and TREA w/o Iso. (the below) to illustrate the impact of the isolation loss  $\mathcal{L}_I$ . Embeddings are projected through t-SNE with Perplexity set to 10 and the Iterations set to 13.)

w/o Iso. in Figure 4. It shows that representations of KG entities in model without the isolation loss are more congested and less distinguishable. It demonstrates the isolation loss can prohibit the clustering of the nodes in KG, which is consistent with the above conjecture.

**Generation.** To examine whether the extraction of the relevant information through the reasoning tree benefits the generation, we conduct the ablation study based on three variants of our complete model, which utilize the whole historical entities, the whole historical utterances and both of the above without extraction, namely *TREA w/o Ent.*, *TREA w/o Utt.*, *TREA w/o EU.* respectively. The results in Table 4 show that deleting either extraction brings a performance decrease on all generation metrics. PPL (Perplexity) is an automatic evaluation metric for the fluency of generations and confidence in the responses. The results of PPL show that the extraction of the relevant information reduced the model confusion. A substantial decrease on Rel. shows that reasoning-guided extraction especially influences the relevance of the generation.

#### 4.7 Evaluation on Long Conversations

We further evaluate TREA in long conversation scenarios. To the best of our knowledge, it is the

Model	Dist-4	Bleu-3	PPL( $\downarrow$ )	Rel.
TREA	0.839	0.013	4.49	2.43
TREA w/o Ent.	0.799	0.012	4.56	2.28
TREA w/o Utt.	0.764	0.011	4.61	2.13
TREA w/o EU.	0.789	0.011	4.78	2.10

Table 4: Evaluation results on the ablation study of the generation task. Fleiss’s kappa values of Rel. all exceed 0.65.

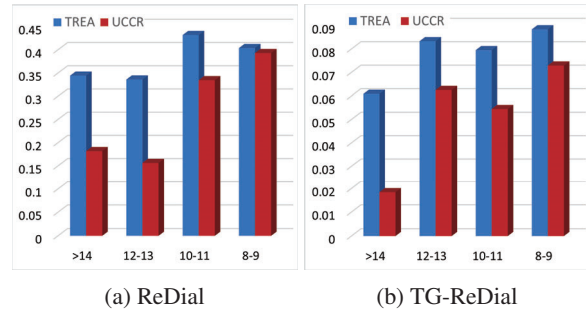


Figure 5: Evaluation results (R@50) of TREA and UCCR on data of different conversation rounds.

first time to discuss this aspect of CRS. When the dialogue becomes longer and more knowledge information appears, if the relationships between knowledge pieces are not clarified, the model is not able to utilize the historical information effectively. We evaluate our TREA and a competitive baseline UCCR on data of different conversation rounds, measured by the metric Recall@50. The results in Figure 5 shows that the performance of UCCR decreases sharply when the conversation rounds exceed 12 in ReDial and 14 in TG-ReDial. On the contrary, the performance of TREA fluctuates less as the number of conversation rounds increases. It indicates that the reasoning process of TREA can illuminate sophisticated relationships between historical entities for a better reference to the current situation, which further proves that nonlinear reasoning with historical hierarchy is vital to modeling user preference, especially when the conversation is long and the informativeness is great.

## 5 Conclusion

In this paper, we propose a novel tree-structure reasoning schema for CRS to clarify the sophisticated relationships between mentioned entities for accurate user modeling. In the constructed reasoning tree, each entity is connected to its cause which motivates the mention of the entity to pro-



vide a clear reference for the current recommendation. The generation module also interacts with the reasoning tree to extract relevant textual information. Extensive experimental results have shown that our approach outperforms several competitive baselines, especially in long conversation scenarios.

## 6 Limitations

The construction of the reasoning tree may be affected by the KG quality since the connection operations are variant with the KG structure. Hence the unsolved problem in Knowledge Graph such as incompleteness or noise could disturb the reasoning process. In the future, we will explore a solution to alleviate the influence of the side information.

## Acknowledgements

This work was supported in part by the National Natural Science Foundation of China under Grant No.62276110, No.62172039 and in part by the fund of Joint Laboratory of HUST and Pingan Property Casualty Research (HPL). The authors would also like to thank the anonymous reviewers for their comments on improving the quality of this paper.

## References

- Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. 2009. Dbpedia - A crystallization point for the web of data. *J. Web Semant.*, 7(3):154–165.
- Qibin Chen, Junyang Lin, Yichang Zhang, Ming Ding, Yukuo Cen, Hongxia Yang, and Jie Tang. 2019. Towards knowledge-based recommender dialog system. pages 1803–1813.
- Zhongxia Chen, Xiting Wang, Xing Xie, Mehul Parsana, Akshay Soni, Xiang Ao, and Enhong Chen. 2020. Towards explainable conversational recommendation. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 2994–3000. ijcai.org.
- Yang Deng, Yaliang Li, Fei Sun, Bolin Ding, and Wai Lam. 2021. Unified conversational recommendation policy learning via graph-based reinforcement learning. pages 1431–1441.
- Yingjie Gu, Xiaoye Qu, Zhefeng Wang, Yi Zheng, Baoxing Huai, and Nicholas Jing Yuan. 2022. Delving deep into regularity: A simple but effective method for chinese named entity recognition. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1863–1873.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.
- Wenqiang Lei, Xiangnan He, Yisong Miao, Qingyun Wu, Richang Hong, Min-Yen Kan, and Tat-Seng Chua. 2020a. Estimation-action-reflection: Towards deep interaction between conversational and recommender systems. In *WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3-7, 2020*, pages 304–312. ACM.
- Wenqiang Lei, Gangyi Zhang, Xiangnan He, Yisong Miao, Xiang Wang, Liang Chen, and Tat-Seng Chua. 2020b. Interactive path reasoning on graph for conversational recommendation. pages 2073–2083.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 110–119. The Association for Computational Linguistics.
- Raymond Li, Samira Ebrahimi Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal. 2018. Towards deep conversational recommendations. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 9748–9758.
- Shuokai Li, Ruobing Xie, Yongchun Zhu, Xiang Ao, Fuzhen Zhuang, and Qing He. 2022. User-centric conversational recommendation with multi-aspect user modeling. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, pages 223–233. ACM.
- Zujie Liang, Huang Hu, Can Xu, Jian Miao, Yingying He, Yining Chen, Xiubo Geng, Fan Liang, and Daxin Jiang. 2021. Learning neural templates for recommender dialogue system. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 7821–7833. Association for Computational Linguistics.
- Zeming Liu, Haifeng Wang, Zhengyu Niu, Hua Wu, and Wanxiang Che. 2021. Durecdial 2.0: A bilingual parallel corpus for conversational recommendation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP*

- 2021, *Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 4335–4347. Association for Computational Linguistics.
- Yu Lu, Junwei Bao, Yan Song, Zichen Ma, Shuguang Cui, Youzheng Wu, and Xiaodong He. 2021. Revcore: Review-augmented conversational recommendation. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pages 1161–1173. Association for Computational Linguistics.
- Wenchang Ma, Ryuichi Takanobu, and Minlie Huang. 2021. Cr-walker: Tree-structured graph reasoning and dialog acts for conversational recommendation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 1839–1851. Association for Computational Linguistics.
- Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. Opendialkg: Explainable conversational reasoning with attention-based walks over knowledge graphs. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 845–854. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pages 311–318. ACL.
- Xiaoye Qu, Yingjie Gu, Qingrong Xia, Zechang Li, Zhefeng Wang, and Baoxing Huai. 2023. A survey on arabic named entity recognition: Past, recent advances, and future trends. *arXiv preprint arXiv:2302.03512*.
- Xiaoye Qu, Jun Zeng, Daizong Liu, Zhefeng Wang, Baoxing Huai, and Pan Zhou. 2022. Distantly-supervised named entity recognition with adaptive teacher learning and fine-grained student ensemble. *arXiv preprint arXiv:2212.06522*.
- Rajdeep Sarkar, Koustava Goswami, Mihael Arcan, and John P. McCrae. 2020. Suggest me a movie for tonight: Leveraging knowledge graphs for conversational recommendation. In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 4179–4189. International Committee on Computational Linguistics.
- Michael Sejr Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In *The Semantic Web - 15th International Conference, ESWC 2018, Heraklion, Crete, Greece, June 3-7, 2018, Proceedings*, volume 10843 of *Lecture Notes in Computer Science*, pages 593–607. Springer.
- Iulian Vlad Serban, Alessandro Sordani, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron C. Courville, and Yoshua Bengio. 2017. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, pages 3295–3301. AAAI Press.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, pages 4444–4451. AAAI Press.
- Yueming Sun and Yi Zhang. 2018. Conversational recommender system. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018*, pages 235–244. ACM.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.
- Xiaolei Wang, Kun Zhou, Ji-Rong Wen, and Wayne Xin Zhao. 2022. Towards unified conversational recommender systems via knowledge-enhanced prompt learning. In *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, pages 1929–1937. ACM.
- Bowen Yang, Cong Han, Yu Li, Lei Zuo, and Zhou Yu. 2022. Improving conversational recommendation systems' quality with context-aware item meta-information. In *Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 38–48. Association for Computational Linguistics.
- Yiming Zhang, Lingfei Wu, Qi Shen, Yitong Pang, Zhihua Wei, Fangli Xu, Bo Long, and Jian Pei. 2022. Multiple choice questions based multi-interest policy learning for conversational recommendation. In *WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022*, pages 2153–2162. ACM.
- Jinfeng Zhou, Bo Wang, Ruifang He, and Yuexian Hou. 2021. CRFR: improving conversational recommender systems via flexible fragments reasoning on knowledge graphs. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 4324–4334. Association for Computational Linguistics.

Kun Zhou, Wayne Xin Zhao, Shuqing Bian, Yuanhang Zhou, Ji-Rong Wen, and Jingsong Yu. 2020a. Improving conversational recommender systems via knowledge graph based semantic fusion. pages 1006–1014.

Kun Zhou, Yuanhang Zhou, Wayne Xin Zhao, Xiaoke Wang, and Ji-Rong Wen. 2020b. Towards topic-guided conversational recommender system. In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 4128–4139. International Committee on Computational Linguistics.

Yuanhang Zhou, Kun Zhou, Wayne Xin Zhao, Cheng Wang, Peng Jiang, and He Hu. 2022. C<sup>2</sup>-crs: Coarse-to-fine contrastive learning for conversational recommender system. In *WSDM '22: The Fifteenth ACM International Conference on Web Search and Data Mining, Virtual Event / Tempe, AZ, USA, February 21 - 25, 2022*, pages 1488–1496. ACM.

## ACL 2023 Responsible NLP Checklist

---

### A For every submission:

- A1. Did you describe the limitations of your work?

6

- A2. Did you discuss any potential risks of your work?

6

- A3. Do the abstract and introduction summarize the paper's main claims?

1

- A4. Have you used AI writing assistants when working on this paper?

*Left blank.*

### B Did you use or create scientific artifacts?

*Left blank.*

- B1. Did you cite the creators of artifacts you used?

*No response.*

- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?

*No response.*

- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?

*No response.*

- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?

*No response.*

- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

*No response.*

- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.

*No response.*

### C Did you run computational experiments?

4

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?

4

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*



- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?  
4
  - C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?  
4
  - C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?  
4
- D  Did you use human annotators (e.g., crowdworkers) or research with human participants?**  
4
- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?  
4
  - D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?  
4
  - D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?  
4
  - D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?  
4
  - D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?  
4