

To Know What User Concerns: Conceptual Knowledge Reasoning for Task-oriented Dialogue Quality Estimation

Anonymous ACL submission

Abstract

Dialogue Quality Estimation (DQE) is crucial in assessing the effects of a conversational consultation system, which has wide applications in E-commerce and Social Media. In task-oriented scenarios, users usually seek personalized consultation about the target subjects they are concerned with rather than general knowledge commonly known by populations. It is essential to identify whether a dialogue solves the user’s questions by task-oriented DQE. Existing studies mainly focus on analyzing dialogue semantics and user sentiment, neglecting to understand what the user is concerned about when requesting a consultation. It may cause fatal errors when the response is emotionally friendly but non-informative. In this paper, we propose a knowledge-enhanced DQE model named **CoReT**, which introduces the **Conceptual Knowledge Reasoning for Task-oriented DQE**. We first design a simple yet efficient entity linking and relation selection module enabling conceptual reasoning from a knowledge graph. Then, we propose a multi-turn textual encoder to capture the contextual information in dialogues. Finally, we introduce a knowledge enhancement module to fuse conceptual reasoning features into contextual embeddings to produce DQE results. For evaluation, we conduct experiments on two real-world datasets in e-commerce consultation systems, the results demonstrate the effectiveness and robustness of **CoReT** compared with the state-of-the-art baselines.

1 Introduction

Dialogue System (DS) is a human-computer interaction technology that aims to simulate natural conversations between humans. It has been widely used in many areas to provide intelligent assistants (Yan et al., 2017), customer service management (Cui et al., 2017), and automated question answering (Zhang et al., 2020). Especially in the e-commerce domain, intelligent assistants play a cru-

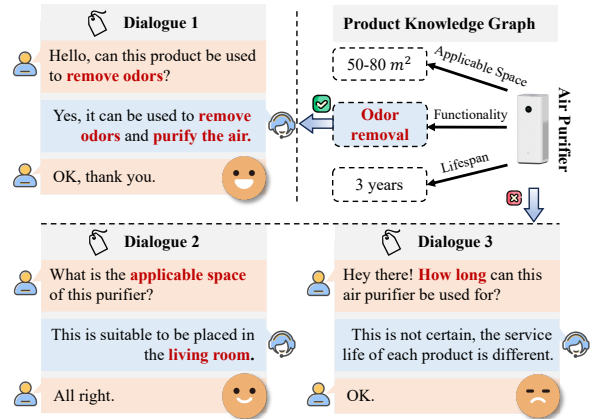


Figure 1: Three dialogues of air purifier. Customers are more satisfied when the server provides accurate answers to product inquiries.

cial role in supporting various customer-oriented services (Li et al., 2017; Ping et al., 2019), such as product consultation, complaints addressing, and feedback collection. The quality of DS is essential for user satisfaction and customer conversion.

However, task-oriented dialogue quality estimation (DQE) remains a challenging problem. It requires not only semantic understanding but also intent detection in textual dialogues, which makes it a complex and difficult task (Fan and Luo, 2020). Recently, task-oriented DQE (Song et al., 2019; Bodigutla et al., 2020; Cai and Chen, 2020) has become a crucial topic in dialogue system research. It has a natural label to measure the quality by user satisfaction. Existing works usually focus on intent detection by modeling user actions in each turn and ultimately fit them to user satisfaction (Sun et al., 2021; Deng et al., 2022b; Kim and Lipani, 2022), or address the semantic understanding by modeling the contextual information within the whole dialogue (Song et al., 2019; Feng et al., 2023). Such work seldom considers an important aspect in task-oriented dialogues, *the subject that the user is inquiring about*, which can provide valuable insights

068 into the specific concerns of the user during the
069 dialogue.

070 Figure 1 presents three dialogue cases in e-
071 commerce consultation systems. The servers who
072 offer accurate and informative answers to user’s
073 queries about the product are more likely to re-
074 ceive better evaluations. The satisfaction of users
075 largely depends on whether the server effectively
076 conveys the necessary information about the sub-
077 ject attributes. This involves the server’s skills in
078 critical information expression during the dialogue.
079 Therefore, utilizing specific knowledge for extract-
080 ing keywords related to subjects in dialogues is a
081 necessary and pivotal work. For example, in the
082 cases of these dialogues, we can capture the key
083 terms in the dialogue by utilizing the triplet \langle *air*
084 *purifier, applicable space, 50-60 square meters* \rangle
085 and \langle *air purifier, lifespan, 2-3 years* \rangle . This in-
086 formation provides a conceptual understanding for
087 servers to infer which aspects of the subjects the
088 customer is concerned about and helps better eval-
089 uate the responses’ quality.

090 Nevertheless, to infer the conceptual knowledge
091 of subjects given a task-oriented dialogue still has
092 to address the following challenges: (1) **Context-**
093 **ual entity matching:** How can we accurately link
094 subject references in dialogues to specific entities
095 in the conceptual knowledge graph? (2) **Differenti-**
096 **ated knowledge reasoning:** How can we perform
097 differentiated knowledge inference for various sub-
098 jects to obtain concepts with different emphases?

099 To address these challenges, we propose **CoReT**
100 to incorporate **C**onceptual knowledge **R**easoning
101 of subjects in dialogues for **T**ask-oriented DQE.
102 Specifically, it consists of three main components,
103 (1) **Conceptual reasoning module:** We trained the
104 TuckER (Balažević et al., 2019) on a large-scale e-
105 commerce knowledge graph OpenBG (Deng et al.,
106 2022a) for conceptual knowledge reasoning. To
107 align inquired subjects with the knowledge graph,
108 we utilize text and semantic similarity matching
109 to link them to corresponding entities. A single
110 subject may be linked to multiple entities, thereby
111 obtaining broader knowledge related to the subject.
112 Additionally, we leverage category-based relations
113 for knowledge inference, enabling differentiated
114 keyword extraction for different types of subjects.
115 By employing knowledge reasoning, we obtain the
116 subjects’ attributes and concepts, which serve as
117 the basis for extracting dialogue keywords. (2) **Hi-**
118 **erarchical Text Mining Module:** To make full use
119 of the information in multi-turn dialogues, we em-

120 ploy a Transformer encoder (Vaswani et al., 2017)
121 to model the semantic actions of both turn-level
122 and dialogue-level. This module enables the model
123 to comprehensively capture the crucial informa-
124 tion within each turn and the historical context
125 across turns. (3) **Knowledge enhancement Mod-**
126 **ule:** Based on the inquired subjects and extracted
127 keywords, we employ parallel multi-head attention
128 mechanisms to enhance the dialogue embeddings.
129 This module aims to integrate conceptual reason-
130 ing features into contextual embeddings, enabling
131 the model to achieve a deeper understanding of the
132 informative and crucial aspects of the dialogue. In
133 summary, our main contributions are as follows:

- 134 • Conceptually, we are the first to model task-
135 oriented DQE by incorporating conceptual
136 knowledge reasoning. Our model captures
137 the attributes of specific subjects and the gen-
138 eralized concepts to make up for the lack of
139 informativeness measurement in DQE.
- 140 • Technically, we propose CoReT, a knowledge-
141 enhanced model that enhances DQE by com-
142 bining conceptual knowledge reasoning with
143 contextual mining.
- 144 • Experimentally, our model achieves the new
145 SOTA on two real-world datasets. Addition-
146 ally, we conduct experiments in multi-task
147 and low-resource scenarios. The results indi-
148 cate that CoReT shows robust and impressive
149 performances in various settings.

150 2 Problem Definition

151 A dialogue can be represented as a sequence of
152 turns $\mathcal{D} = \{(u_{u_1}, u_{s_1}), \dots, (u_{u_T}, u_{s_T})\}$, where
153 (u_{u_t}, u_{s_t}) represents the user’s question and the
154 server’s answer in the t -th turn. A complete dia-
155 logue usually starts with a question or consultation
156 about a specific subject. Task-oriented DQE re-
157 quires not only understanding the informativeness
158 of the response whether it addresses the question
159 but also the sentimental reactions to customers’
160 concerns. Therefore, the evaluation of the dialogue
161 quality about how it meets customer satisfaction
162 is summarized as follows: given a dialogue \mathcal{D} , the
163 subject information \mathcal{P} , and the conceptual knowl-
164 edge graph \mathcal{G} , the objective is to model an estimator
165 \mathcal{E} that accurately predicts the user satisfaction \mathcal{Y}^s
166 by the end of the dialogue session. This process

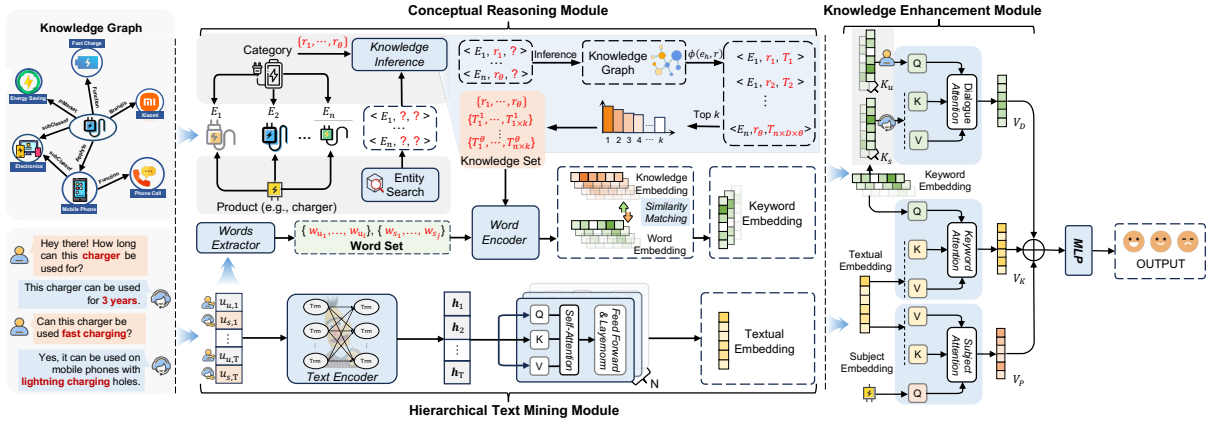


Figure 2: Overview of CoReT: (1) Conceptual Reasoning Module perform subject knowledge inference and extract the keywords from the dialogue; (2) Hierarchical Text Mining Module encodes the dialogue; (3) Knowledge Enhancement Module fuse the keywords, contextual representation and subject information.

can be formalized as:

$$y^s = \mathcal{E}(\mathcal{D}, \mathcal{P}, \mathcal{G}) \quad (1)$$

3 Methods

Figure 2 shows the architecture of our proposed model. It is composed of three parts: (1) **Conceptual Reasoning module**: We begin by pretraining a knowledge-inferring model on the knowledge graph. We link each subject to its corresponding entities as head entities to perform knowledge inference. Using the head entities and relations, we predict the tail entities on the inferring model, thereby obtaining key attributes and concepts for each subject. This enables us to extract the keywords from the dialogue. (2) **Hierarchical Text Mining Module**: Utilizing the utterance encoder and Transformer for turn-level and dialogue-level encoding to obtain dialogue contextual representation. (3) **Knowledge Enhancement module**: Applying attention mechanisms to fuse different representations with distinct focuses.

3.1 Conceptual Reasoning Module

To conduct unique knowledge inference for each subject, thus obtaining its attributes and concepts, we introduce the Conceptual Reasoning Module. It enables the model to effectively capture informative keywords in a dialogue.

3.1.1 Pretraining Knowledge Inferring Model

We select the TuckER (Balažević et al., 2019) model for knowledge graph inference. For a triple (e_h, r, e_t) in the knowledge graph, we define the embedding vectors for the head entity and relation as $\mathbf{v}_h \in \mathbb{R}^{d_e}$ and $\mathbf{v}_r \in \mathbb{R}^{d_r}$, respectively. We also

introduce the parameter matrix $\mathbf{W} \in \mathbb{R}^{d_e \times d_e \times d_r}$ and the unified representation matrix $\mathbf{U} \in \mathbb{R}^{d_e \times M}$ that contains embeddings for all entities, where M denotes the number of entities. Given e_h and r , the inferring of tail entity is as follows:

$$\phi(e_h, r) = \mathbf{W} \times \mathbf{v}_r \times \mathbf{v}_h \times \mathbf{U} \quad (2)$$

Here, $\phi(e_h, r)$ represents the probability distribution of the predicted entities according to the head entity and relation: $(e_h, r, ?)$. We consider the entity with the highest probability as the predicted result and use the cross-entropy loss function to train the knowledge inferring model. In practical applications, we choose the top- K entities with the highest probabilities as potential tail entities. For each subject, we consider both the relations $\{r\}$ utilized in the inference and the resulting tail entities $\{T\}$ as the key attributes $\{attr_1, \dots, attr_h\}$.

3.1.2 Matching Subject Entities & Selecting Relations

Before inferring the conceptual knowledge corresponding to the inquiry subjects, we have to map the subjects to the entity nodes in knowledge graphs. Considering the efficiency of subject entity matching between dialogues and knowledge bases, we propose a simple yet efficient method to match the subjects \mathcal{P} mentioned in the dialogue to the corresponding entity E in the knowledge graph by judging whether their text is the same ($\mathcal{P} = E$) or their semantic cosine similarity is greater than a threshold ($\cos_smi(\text{BERT}(\mathcal{P}), \text{BERT}(E)) > \tau$). To fully explore the relevant entities of a subject, we link each subject to multiple entities $\{E_1, \dots, E_n\}$, using them as candidate head entities in the inference process. This approach allows us to

leverage a broader range of knowledge and capture a comprehensive understanding of the subject.

To acquire the relations utilized during the inference process for different subjects, we design a statistical-based heuristic algorithm. For different categories of subjects, it retrieves and counts the relations that appear in the training set dialogues. Based on the frequency of occurrence, we ultimately select the top- θ most frequently occurring relations within each category.

3.1.3 Extracting Keywords in Dialogues

Based on the candidate entities matched to the subjects and corresponding relation sets describing the inferring aspects of differentiated subject categories, we utilize the key attributes $\{attr_1, \dots, attr_h\}$ which contains both the relations and inferred tail entities to extract keywords in the dialogue. Then, we segment the dialogue utterances to obtain the word sequences for the user and the server, denoted as $\{w_{u_1}, \dots, w_{u_i}\}$ and $\{w_{s_1}, \dots, w_{s_j}\}$ respectively. A semantic matching mechanism is designed to calculate the similarity between the segmented words $\{w\}$ and the entity attributes $\{attr\}$ using the following formula:

$$\text{Sim}(w, attr) = \cos_sim(\text{BERT}(w), \text{BERT}(attr)) \quad (3)$$

The BERT model is employed to encode the words. If the similarity between a $(w, attr)$ pair exceeds a threshold τ , we consider the word w as a keyword. As a result, for each dialogue, we have the keywords $\{w_{uk_1}, \dots, w_{uk_n}\}$ and $\{w_{sk_1}, \dots, w_{sk_m}\}$ from user and server respectively.

3.2 Hierarchical Text Mining Module

To address multi-turn customer service dialogues, we obtain contextual representations that capture both the turn-level and the dialogue-level information in this module. Firstly, for turn-level context encoding, we utilize the BERT model as the text encoder. The user and the server utterances (u_{ut}, u_{st}) are concatenated as a complete sequence with a special token $[SEP]$ to feed the text encoder. Then the turn-level representation is calculated as follows:

$$\begin{aligned} \mathbf{h}_t &= \text{TextEncoder}(u_{ut}, u_{st}) \\ &= \text{BERT}([CLS]u_{ut}[SEP]u_{st}[SEP]) \end{aligned} \quad (4)$$

By calculating all \mathbf{h}_t for each turn, we can obtain a set of turn-level representations $\mathbf{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_T\}$, where T is the number of turns in the dialogue.

To capture the inter-turn relationships and obtain an overall representation of the dialogue, we employ a dialogue-level encoder with L transformer layers. Similar to the turn-level encoder, we also add position embedding vectors to ensure the temporal information between different turns. The inputs are constructed as follows:

$$\mathbf{H}^{(0)} = \{\mathbf{h}_t + \text{position_encoding}(t) | 1 \leq t \leq T\} \quad (5)$$

where the position encoding is computed using sine and cosine functions following the approach in (Vaswani et al., 2017). Then we utilize the Multi-Head Attention mechanism (MHA) to fuse information across different turns.

The computation of the dialogue-level encoder that we employ is as follows:

$$\begin{aligned} \mathbf{H}' &= \text{MHA}(\mathbf{H}^{(l)}, \mathbf{H}^{(l)}, \mathbf{H}^{(l)}) \\ \mathbf{H}^{(l+1)} &= \text{FFN}(\mathbf{H}' + \mathbf{H}^{(l)}) + \mathbf{H}' + \mathbf{H}^{(l)} \end{aligned} \quad (6)$$

where l represents the l -th encoder layer. The Feed Forward Network (FFN) is also utilized to facilitate information flow between the encoder layers. Within the FFN, we use the ReLU function as the activation function between the two linear layers.

Finally, we retrieve the last layer's output of the encoder, denoted as $\mathbf{H}^{(L)} = \{\mathbf{h}_0^{(L)}, \dots, \mathbf{h}_T^{(L)}\}$. From $\mathbf{H}^{(L)}$, we select the encoded information of the last turn $\mathbf{h}_T^{(L)}$ as the contextual representation of the whole dialogue, denoted as \mathbf{D}_e .

3.3 Knowledge Enhancement Module

In this module, we employ multiple multi-head attention blocks to effectively integrate keyword information, dialogue text information, and subject information. This integration allows for comprehensive exploration of key information within the dialogue from various perspectives.

Firstly, we introduce the keyword-level feature fusion. Given the obtained keywords $\{w_{uk_1}, \dots, w_{uk_n}\}$ and $\{w_{sk_1}, \dots, w_{sk_m}\}$ from the conceptual reasoning module, we utilize the BERT model to capture the conceptual features about subject keywords. The keyword representations are calculated as follows:

$$\begin{aligned} \mathbf{K}_u &= \text{BERT}([CLS]w_{uk_1}[SEP] \dots w_{uk_n}[SEP]) \\ \mathbf{K}_s &= \text{BERT}([CLS]w_{sk_1}[SEP] \dots w_{sk_m}[SEP]) \end{aligned} \quad (7)$$

Then we apply a multi-head attention block to

320 fuse the two vectors \mathbf{K}_u and \mathbf{K}_s :

$$321 \quad \begin{aligned} \mathbf{V}_D &= \text{DialogueAttention}(\mathbf{K}_u, \mathbf{K}_s, \mathbf{K}_s) \\ &= \text{MHA}(\mathbf{K}_u, \mathbf{K}_s, \mathbf{K}_s) \end{aligned} \quad (8)$$

322 where we obtain \mathbf{V}_D to provide essential interac-
323 tive information between the user and the server in
324 a dialogue.

325 Furthermore, we specifically conduct informa-
326 tion mining to target contextual representation with
327 integrated keyword representation:

$$328 \quad \begin{aligned} \mathbf{V}_K &= \text{KeywordAttention}(\mathbf{K}_{us}, \mathbf{D}_e, \mathbf{D}_e) \\ &= \text{MHA}(\mathbf{K}_{us}, \mathbf{D}_e, \mathbf{D}_e) \end{aligned} \quad (9)$$

329 where \mathbf{K}_{us} is the integrated keyword representation
330 by averaging \mathbf{K}_u and \mathbf{K}_s .

331 Finally, an embedding layer is utilized to rep-
332 resent each unique subject \mathbf{P}_e . We then fuse the
333 subject embedding with the contextual representa-
334 tion as follows:

$$335 \quad \begin{aligned} \mathbf{V}_P &= \text{SubjectAttention}(\mathbf{P}_e, \mathbf{D}_e, \mathbf{D}_e) \\ &= \text{MHA}(\mathbf{P}_e, \mathbf{D}_e, \mathbf{D}_e) \end{aligned} \quad (10)$$

336 In summary, we obtain three representations in
337 this module: \mathbf{V}_D , \mathbf{V}_K , and \mathbf{V}_P . Among them, \mathbf{V}_D
338 represents the key interaction between the user and
339 the server. \mathbf{V}_K represents the contextual represen-
340 tation enhanced by the subject-related keywords.
341 And for \mathbf{V}_P , it represents the contextual represen-
342 tation autonomously enhanced by the subject in-
343 formation. These three types of representations
344 enable us to explore the dialogue from multiple
345 perspectives.

346 3.4 Learning Objectives

347 In the output module, we concatenate the repre-
348 sentations \mathbf{V}_D , \mathbf{V}_K , and \mathbf{V}_P and feed them into a
349 multi-layer perceptron (MLP) for the final satisfac-
350 tion estimation:

$$351 \quad y_p^s = \text{softmax}(\text{MLP}(\text{Concat}(\mathbf{V}_D, \mathbf{V}_K, \mathbf{V}_P))) \quad (11)$$

352 where y_p^s represents the prediction results produced
353 by our model. We utilize the cross-entropy loss as
354 the objective function to train the proposed model:

$$355 \quad \mathcal{L} = -y^s \log(y_p^s) \quad (12)$$

356 where y^s represents the ground truth. During the
357 training process, the parameters of the pre-trained
358 BERT model are also fine-tuned.

Table 1: The statistics of datasets.

Dataset	# Dialogues	# Aver. utterance	# Labels
JDDC	3,300	10.9	3
EDP	8,115	10.2	5

359 4 Experiments

360 4.1 Dataset & Evaluation metrics

361 To evaluate the effectiveness of our model, we con-
362 duct the experiments on an open-source dataset
363 (**JDDC** (Chen et al., 2019)) and a real-world cus-
364 tomer service dialogue dataset (**EDP**) collected
365 from 300 users as volunteers. The dialogues are
366 conversational consultations in e-commerce plat-
367 forms. We use the user satisfaction score as the
368 dialogue quality label for evaluations. The statis-
369 tics of the datasets are shown in Table 1.

370 We consider the inquired products in the dia-
371 logue as the subjects, which cover a total of 602
372 distinct products in the EDP dataset. When there is
373 no specific subject in a dialogue, we general tokens
374 of the most frequent products within each category
375 to serve as the possible subjects. For the conceptual
376 reasoning, we utilize OpenBG (Deng et al., 2022a)
377 as the product knowledge graph.

378 We adopted accuracy, macro-average precision,
379 recall, and F1-score as evaluation metrics in our
380 experiments, which are consistent with previous
381 works (Ye et al., 2023; Feng et al., 2023; Deng
382 et al., 2022b). The implementation code of our
383 model is available here¹.

384 4.2 Baselines

385 We selected the following models as baselines:

386 **HAN** (Yang et al., 2016) It utilizes a two-level
387 attention mechanism to capture information at dif-
388 ferent granularities in dialogues.

389 **Transformer** (Vaswani et al., 2017) It passes the
390 dialogue embedding as input to the Transformer
391 for representation learning.

392 **BERT** (Devlin et al., 2018) It embeds multi-turn
393 dialogue by concatenating each turn.

394 **Speaker** (He et al., 2021) It leverages the mod-
395 eling of interactions between the participants to
396 recognize the dialogue acts.

397 **CDCN** (Li et al., 2020b) It employs a dynamic
398 convolutional network as an utterance encoder to
399 capture local information in dialogues.

¹<https://anonymous.4open.science/r/CoRe-USE-582C/>

400 **SG-USM** (Feng et al., 2023) It adopts an atten- 447
 401 tion mechanism to model task fulfillment in task-
 402 oriented dialogues.
 403 **USDA** (Deng et al., 2022b) It employs a hierarchi- 448
 404 cal Transformer structure for dialogue encoding. 449
 405 **ASAP** (Ye et al., 2023) It applies the Hawkes pro- 450
 406 cess to model the intent changes in DQE tasks. 451
 407 To ensure fairness, we employ the BERT model 452
 408 as the backbone encoder for all baseline models. 453

409 4.3 Parameter Settings 454

410 The base BERT consists of 12 Transformer lay- 455
 411 ers and outputs a final dimensionality of 768. In 456
 412 the conceptual reasoning module, we set the co- 457
 413 sine similarity threshold τ to 0.80. Each subject 458
 414 is mapped to 10 different entities and we select 459
 415 50 relations for each subject category. For each 460
 416 inference, we select the top 5 tail entities as the 461
 417 reliable results. To verify the reliability of the rea- 462
 418 soning model, we conducted tail entity prediction 463
 419 tasks on the benchmark of OpenBG. The model 464
 420 achieved impressive results with hits@10=71.0%, 465
 421 hits@3=59.0%, and hits@1=41.4%. For the 466
 422 utterances-level representation module, we set the 467
 423 number of Transformer layers to 3 and the number 468
 424 of attention heads to 8. The dropout probability of 469
 425 MLP is set to 0.1. We utilized the Adam optimizer 470
 426 and trained the model for 10 epochs. 471

427 4.4 Overall Performance 472

428 The results of DQE on the EDP and JDDC datasets 473
 429 are shown in Table 2. It can be observed that our 474
 430 model significantly and consistently outperforms 475
 431 all baselines on both datasets. Since JDDC lacks 476
 432 explicit subject information for each dialogue, we 477
 433 utilize the most frequent products within each cat- 478
 434 egory as the head entities used in the inference 479
 435 process for experiments conducted on this dataset. 480
 436 Even without explicit subject information, our ap- 481
 437 proach still achieves a 1.5% improvement in F1- 482
 438 score compared to ASAP, which is the state-of- 483
 439 the-art model. Moreover, for EDP where subject 484
 440 information is available, our approach consistently 485
 441 outperforms other baseline models with an average 486
 442 10% higher F1-score. This indicates that our mod- 487
 443 eling of subject information enables the model to 488
 444 focus on the crucial aspects that impact user satis- 489
 445 faction in dialogues, thereby effectively enhancing 490
 446 the accuracy of dialogue quality evaluation. 491

4.5 Ablation Study 447

In the ablation experiments, we conducted a total of 448
 five groups of experiments. Among them, “w/o text 449
 mining” refers to the removal of dialogue-level en- 450
 coding from the Hierarchical Text Mining Module, 451
 and directly using the [CLS] token output by the 452
 BERT model for dialogue text representation. The 453
 results of ablation experiments on two datasets are 454
 shown in Figure 4. It can be observed that utilizing 455
 the transformer model for dialogue-level encoding 456
 is beneficial for the overall DQE. Furthermore, we 457
 also conducted ablation experiments on each of 458
 the three attention modules within the Knowledge 459
 Enhancement Module. The experimental results 460
 also confirmed that all three attention modules we 461
 designed positively contribute to the experimental 462
 outcomes. The final group of ablation experiments 463
 targeted the Conceptual Reasoning Module. In this 464
 experiment, we performed inference on the head 465
 entities using all relations in the knowledge graph, 466
 rather than solely relying on the most relevant and 467
 important relations for a particular subject category. 468
 From the experimental results, it is evident that if 469
 the relationships in the knowledge graph are not 470
 filtered, the knowledge reasoning module may in- 471
 troduce additional noise, thus negatively impacting 472
 the experimental results. 473

4.6 Robustness Assessment 474

4.6.1 Low-resource Scenarios 475

Due to the scarcity of high-quality training data for 476
 customer service dialogues in the real world, we 477
 also conducted evaluations of model performance 478
 in low-resource scenarios. In this experiment, we 479
 select USDA and ASAP, which performed well 480
 in the DQE task, as the baselines. We limited the 481
 number of dialogues used for training to 500, 1,000, 482
 2,000, and 4,000, respectively, and conducted mul- 483
 tiple experiments by randomly sampling from the 484
 two datasets for 5 times to test the stability of model 485
 performance. The results in the low-resource sce- 486
 nario are shown in Figure 3. The lines in the 487
 graph represent the average performance of the 488
 models under different sampling quantities, while 489
 the shaded areas represent the upper and lower 490
 ranges of model performance under different sam- 491
 plings. From the graph, it can be observed that 492
 our model performs better in terms of stability and 493
 performance by incorporating subject knowledge. 494

Table 2: Comparison of overall experimental results between CoReT and baselines. The **bold** demonstrates the best performance, while the second best performance is indicated with an underline.

Model	EDP				JDDC			
	Accuracy	Precision	Recall	F1-score	Accuracy	Precision	Recall	F1-score
HAN	56.3	39.3	39.9	38.2	58.4	54.2	50.1	52.0
Transformer	61.7	57.2	50.1	51.5	58.1	58.8	64.9	58.9
BERT	63.9	55.2	53.7	53.8	60.4	59.8	58.8	59.5
Speaker	59.7	55.7	47.9	50.0	63.4	61.8	62.2	61.9
CDCN	60.4	55.4	49.1	48.3	62.4	59.1	56.1	57.2
SG-USM	61.8	49.4	51.2	50.1	63.3	63.1	64.1	63.5
USDA	62.9	56.1	52.9	54.0	61.8	<u>62.8</u>	63.7	61.7
ASAP	<u>66.8</u>	<u>61.4</u>	<u>54.4</u>	<u>56.0</u>	<u>64.9</u>	<u>62.3</u>	<u>65.4</u>	<u>63.5</u>
CoReT	71.3	62.5	61.4	60.5	65.4	63.7	67.3	65.0

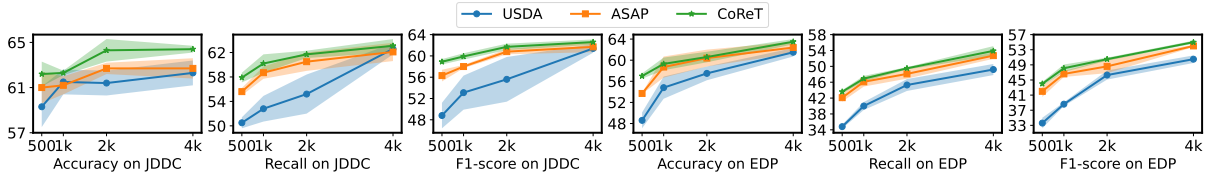


Figure 3: Experimental results in low-resource scenarios.

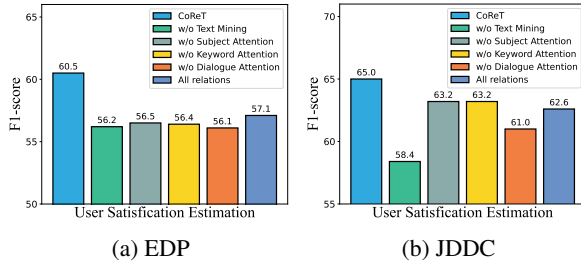


Figure 4: Ablation study on two datasets.

Table 3: Comparison of multi-task experimental results.

Model	DQE			DAR		
	Acc	Recall	F1	Acc	Recall	F1
JointDAS	58.5	55.1	55.4	63.4	43.6	41.1
Co-GAT	60.6	63.7	61.0	66.7	48.9	47.5
JointUSE	63.8	58.6	59.2	66.8	48.7	47.3
Speaker	<u>64.8</u>	60.1	60.2	65.8	47.0	45.8
USDA	63.0	65.7	<u>62.6</u>	<u>69.7</u>	<u>53.0</u>	<u>51.3</u>
ASAP	64.0	61.8	61.7	68.2	50.3	48.5
CoReT	65.2	<u>64.7</u>	63.8	70.1	53.2	51.5

4.6.2 Multi-task Scenarios

To evaluate the practicality of our model in real-world scenarios and enable the platform to track the real-time changes in user demands, we test our model on the JDDC dataset for the Dialogue Act Recognition (DAR) task. The DAR task aims to classify the user’s intent for each turn in a dialogue. In our experiment, we conducted joint training of our model on both DQE and DAR tasks to obtain experimental results. To ensure fair performance comparison among different models, we selected three baseline models, namely JointDAS (Cerisara et al., 2018), Co-GAT (Qin et al., 2021), and JointUSE (Bodigutla et al., 2020), which are designed for multi-task learning. The experimental results in the multitasking scenario are shown in Table 3. Our model outperformed the other baselines in both tasks, indicating its excellent adaptability in turn-level user intent recognition tasks.

Table 4: Results Comparison of LLM and KG.

Knowledge Source	Acc	Recall	F1	Time
GPT-4	66.3	55.9	57.4	1.74s
KG	71.3	61.4	60.5	0.53ms

4.7 Discussion on LLM vs. KG

To investigate how well the LLMs perform in conceptual reasoning rather than a knowledge graph, we use GPT-4 to replace the knowledge graph(KG) for reasoning on subjects. The head entity and the corresponding relations are used to construct prompts of LLMs. The experimental results on the EDP are shown in Table 4. It can be seen that the knowledge obtained from KG is more conducive to accurately estimate the dialogue quality. Moreover, with cuda acceleration, each reasoning on the KG only takes 0.53ms, which is much faster than using LLM inferring APIs.

4.8 Case Study

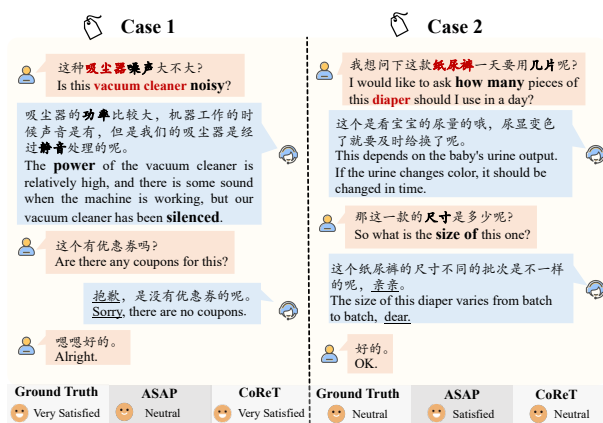


Figure 5: Two dialogue cases selected from the dataset. The mentioned products in the dialogue are highlighted in red, the product-related keywords are highlighted in bold, and the sentiment-related words are indicated with an underline.

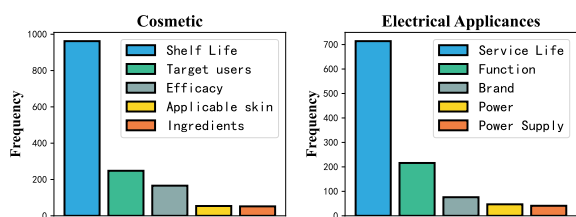


Figure 6: Relation selection of two product categories.

To further illustrate the working principle of our model in dialogue quality evaluation, We selected two dialogue examples from the dataset which are presented in Figure 5. It can be observed that models that rely solely on text analysis tend to focus more on emotion-related words in the dialogue. If the customer service attitude is friendly, it may receive a higher satisfaction rating, while neglecting whether the server answered the user's inquiry about products. However, our model can effectively address this issue and provide accurate evaluations.

We further select two examples to showcase the relations obtained by the algorithm for different product categories, as shown in Figure 6. Our algorithm can provide different relations based on different product categories and the selected relations are closely related to the corresponding subjects.

5 Related Work

5.1 Dialogue Quality Estimation

As intelligent dialogue systems are being widely applied in service platforms, Dialogue Quality Estimation is receiving more and more attention. Due

to the high-cost and time-consuming characteristics of manual evaluation, recent works have increasingly relied on deep learning methods. Such methods follow two main lines. Some researchers focus on modeling user intents at the turn level to assess the dialogue quality (Kim and Lipani, 2022; Deng et al., 2022b). The other methods analyze at the dialogue level, focusing on contextual information the dialogue (Mendonca et al., 2022; Gupta et al., 2021; Ye et al., 2023). However, current methods seldom consider effectively modeling the subjects inquired about in task-oriented dialogues, which can imply the specific concerns of users. This limitation hinders the models' ability to determine whether the dialogue quality meets users' satisfaction.

5.2 Knowledge-enhanced Dialogue Understanding

In the field of natural language processing, knowledge is increasingly being used in various works to enhance the model's comprehension. Some works proposed to encode knowledge by pre-training and then apply the embedding to downstream tasks, such as Know-BERT (Peters et al., 2019b), KEPLER (Wang et al., 2021), etc. While some other studies directly use discrete triplets for knowledge enhancement. This type of work first retrieves relevant knowledge from the knowledge graph and then applies it to the model (Peters et al., 2019a; Li et al., 2020a; Gao et al., 2019). However, existing works on knowledge enhancement mostly focus on utilizing external knowledge to provide additional information, neglecting the exploration of the key content inherent in data.

6 Conclusion

In this study, we propose CoReT, which incorporates conceptual knowledge reasoning of subjects into dialogue quality estimation in multi-turn e-commerce dialogues. By capturing product-related concepts and attributes, our model can focus on the key content related to subject inquiries in the dialogue, which is also the part that concerns users the most. We conducted experiments on the EDP and JDDC datasets. The results demonstrated that our model achieved state-of-art performance on both datasets. To validate the impact of the knowledge on model robustness, we also conduct robustness tests in multi-task and low-resource scenarios. The results demonstrate that our model exhibits stable and reliable performance across various scenarios.

599 Limitations

600 While our model achieves the new state-of-the-art
601 performance, it still has several limitations. Firstly,
602 although we employ the hierarchical structure to
603 model multi-turn dialogues, the abstractive extrac-
604 tion of key sentences in such dialogues still leaves
605 an open issue when faced with long texts. Secondly,
606 when there are no explicit subjects attached to given
607 dialogues, we proposed to use a set of general sub-
608 jects that frequently occurred in all dialogues and
609 provide the model with commonsense reasoning
610 comprehension corresponding to the e-commerce
611 scenarios. Ideally, the framework should be able
612 to understand what subjects the user is concerned
613 about automatically from the whole contextual in-
614 formation rather than certain words at the begin-
615 ning of the dialogues. Lastly, the evaluations are
616 conducted in e-commerce consultation scenarios,
617 the scalability of our model could be further ana-
618 lyzed in other domains.

619 Ethics Statement

620 In this work, we do not use any persona profiles
621 or other user information for modeling. We do not
622 collect or handle any personal data in our exper-
623 iments. When applying our model to real-world
624 applications, there is no need to capture any other
625 information except the raw dialogues. The imple-
626 mentation code of our model is publicly available
627 and compliant with anonymity standards.

628 References

629 Ivana Balažević, Carl Allen, and Timothy M
630 Hospedales. 2019. Tucker: Tensor factorization
631 for knowledge graph completion. *arXiv preprint*
632 *arXiv:1901.09590*.

633 Praveen Kumar Bodigutla, Aditya Tiwari, Josep Valls
634 Vargas, Lazaros Polymenakos, and Spyros Mat-
635 soukas. 2020. Joint turn and dialogue level user sat-
636 isfaction estimation on multi-domain conversations.
637 *arXiv preprint arXiv:2010.02495*.

638 Wanling Cai and Li Chen. 2020. Predicting user intents
639 and satisfaction with dialogue-based conversational
640 recommendations. In *Proceedings of the 28th ACM*
641 *Conference on User Modeling, Adaptation and Per-*
642 *sonalization*, pages 33–42.

643 Christophe Cerisara, Somayeh Jafaritazehjani, Ade-
644 dayo Oluokun, and Hoa Le. 2018. Multi-task dialog
645 act and sentiment recognition on mastodon. *arXiv*
646 *preprint arXiv:1807.05013*.

Meng Chen, Ruixue Liu, Lei Shen, Shaozu Yuan,
Jingyan Zhou, Youzheng Wu, Xiaodong He, and
Bowen Zhou. 2019. The jddc corpus: A large-scale
multi-turn chinese dialogue dataset for e-commerce
customer service. *arXiv preprint arXiv:1911.09969*.

Lei Cui, Shaohan Huang, Furu Wei, Chuanqi Tan, Chao-
qun Duan, and Ming Zhou. 2017. Superagent: A
customer service chatbot for e-commerce websites.
In *Proceedings of ACL 2017, system demonstrations*,
pages 97–102.

Shumin Deng, Hui Chen, Zhoubo Li, Feiyu Xiong,
Qiang Chen, Mosha Chen, Xiangwen Liu, Jiaoyan
Chen, Jeff Z Pan, Huajun Chen, et al. 2022a. Con-
struction and applications of open business knowl-
edge graph. *arXiv preprint arXiv:2209.15214*.

Yang Deng, Wenxuan Zhang, Wai Lam, Hong Cheng,
and Helen Meng. 2022b. User satisfaction estima-
tion with sequential dialogue act modeling in goal-
oriented conversational systems. In *Proceedings of*
the ACM Web Conference 2022, pages 2998–3008.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
Kristina Toutanova. 2018. Bert: Pre-training of deep
bidirectional transformers for language understand-
ing. *arXiv preprint arXiv:1810.04805*.

Yifan Fan and Xudong Luo. 2020. A survey of dialogue
system evaluation. In *2020 IEEE 32nd International*
Conference on Tools with Artificial Intelligence (IC-
TAI), pages 1202–1209. IEEE.

Yue Feng, Yunlong Jiao, Animesh Prasad, Niko-
laos Aletras, Emine Yilmaz, and Gabriella Kazai.
2023. Schema-guided user satisfaction model-
ing for task-oriented dialogues. *arXiv preprint*
arXiv:2305.16798.

Shen Gao, Zhaochun Ren, Yihong Zhao, Dongyan Zhao,
Dawei Yin, and Rui Yan. 2019. Product-aware an-
swer generation in e-commerce question-answering.
In *Proceedings of the Twelfth ACM International*
Conference on Web Search and Data Mining, pages
429–437.

Saurabh Gupta, Xing Fan, Derek Liu, Benjamin Yao,
Yuan Ling, Kun Zhou, Tuan-Hung Pham, and Chen-
lei Edward Guo. 2021. Robertaiq: An efficient frame-
work for automatic interaction quality estimation of
dialogue systems.

Zihao He, Leili Tavabi, Kristina Lerman, and Mo-
hammad Soleymani. 2021. Speaker turn model-
ing for dialogue act classification. *arXiv preprint*
arXiv:2109.05056.

To Eun Kim and Aldo Lipani. 2022. A multi-task based
neural model to simulate users in goal oriented di-
alogue systems. In *Proceedings of the 45th Inter-*
national ACM SIGIR Conference on Research and
Development in Information Retrieval, pages 2115–
2119.

701	Feng-Lin Li, Hehong Chen, Guohai Xu, Tian Qiu, Feng	Weiwei Sun, Shuo Zhang, Krisztian Balog, Zhaochun	759
702	Ji, Ji Zhang, and Haiqing Chen. 2020a. Alimekg: Do-	Ren, Pengjie Ren, Zhumin Chen, and Maarten de Ri-	760
703	main knowledge graph construction and application	ijke. 2021. Simulating user satisfaction for the evalu-	761
704	in e-commerce. In <i>Proceedings of the 29th ACM In-</i>	ation of task-oriented dialogue systems. In <i>Proceed-</i>	762
705	<i>ternational Conference on Information & Knowledge</i>	<i>ings of the 44th International ACM SIGIR Confer-</i>	763
706	<i>Management</i> , pages 2581–2588.	<i>ence on Research and Development in Information</i>	764
		<i>Retrieval</i> , pages 2499–2506.	765
707	Feng-Lin Li, Minghui Qiu, Haiqing Chen, Xiongwei	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	766
708	Wang, Xing Gao, Jun Huang, Juwei Ren, Zhongzhou	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	767
709	Zhao, Weipeng Zhao, Lei Wang, et al. 2017. Alime	Kaiser, and Illia Polosukhin. 2017. Attention is all	768
710	assist: An intelligent assistant for creating an inno-	you need. <i>Advances in neural information processing</i>	769
711	vative e-commerce experience. In <i>Proceedings of</i>	<i>systems</i> , 30.	770
712	<i>the 2017 ACM on Conference on Information and</i>		
713	<i>Knowledge Management</i> , pages 2495–2498.		
714	Jingye Li, Hao Fei, and Donghong Ji. 2020b. Model-	Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan	771
715	ing local contexts for joint dialogue act recognition	Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. 2021.	772
716	and sentiment classification with bi-channel dynamic	Kepler: A unified model for knowledge embedding	773
717	convolutions. In <i>Proceedings of the 28th International</i>	and pre-trained language representation. <i>Transac-</i>	774
718	<i>Conference on Computational Linguistics</i> , pages 616–	<i>tions of the Association for Computational Linguis-</i>	775
719	626.	<i>tics</i> , 9:176–194.	776
720	John Mendonca, Alon Lavie, and Isabel Trancoso. 2022.	Zhao Yan, Nan Duan, Peng Chen, Ming Zhou, Jianshe	777
721	Qualityadapt: an automatic dialogue quality estima-	Zhou, and Zhoujun Li. 2017. Building task-oriented	778
722	tion framework. In <i>Proceedings of the 23rd Annual</i>	dialogue systems for online shopping. In <i>Proceed-</i>	779
723	<i>Meeting of the Special Interest Group on Discourse</i>	<i>ings of the AAAI Conference on Artificial Intelligence</i> ,	780
724	<i>and Dialogue</i> , pages 83–90.	volume 31.	781
725	Matthew E. Peters, Mark Neumann, Robert Logan, Roy	Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He,	782
726	Schwartz, Vidur Joshi, Sameer Singh, and Noah A.	Alex Smola, and Eduard Hovy. 2016. Hierarchical at-	783
727	Smith. 2019a. Knowledge enhanced contextual word	tention networks for document classification. In <i>Pro-</i>	784
728	representations. In <i>Proceedings of the 2019 Confer-</i>	<i>ceedings of the 2016 conference of the North Ameri-</i>	785
729	<i>ence on Empirical Methods in Natural Language Pro-</i>	<i>can chapter of the association for computational lin-</i>	786
730	<i>cessing and the 9th International Joint Conference</i>	<i>guistics: human language technologies</i> , pages 1480–	787
731	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	1489.	788
732	pages 43–54, Hong Kong, China. Association for	Fanghua Ye, Zhiyuan Hu, and Emine Yilmaz. 2023.	789
733	Computational Linguistics.	Modeling user satisfaction dynamics in dialogue via	790
		hawkes process. <i>arXiv preprint arXiv:2305.12594</i> .	791
734	Matthew E Peters, Mark Neumann, Robert L Lo-	Weisheng Zhang, Kaisong Song, Yangyang Kang,	792
735	gan IV, Roy Schwartz, Vidur Joshi, Sameer Singh,	Zhongqing Wang, Changlong Sun, Xiaozhong Liu,	793
736	and Noah A Smith. 2019b. Knowledge enhanced	Shoushan Li, Min Zhang, and Luo Si. 2020. Multi-	794
737	contextual word representations. <i>arXiv preprint</i>	turn dialogue generation in e-commerce platform	795
738	<i>arXiv:1909.04164</i> .	with the context of historical dialogue. In <i>Findings</i>	796
739	Ng Lian Ping et al. 2019. Constructs for artificial in-	<i>of the Association for Computational Linguistics:</i>	797
740	telligence customer service in e-commerce. In <i>2019</i>	<i>EMNLP 2020</i> , pages 1981–1990.	798
741	<i>6th International Conference on Research and Inno-</i>		
742	<i>vation in Information Systems (ICRIIS)</i> , pages 1–6.		
743	IEEE.		
744	Libo Qin, Zhouyang Li, Wanxiang Che, Minheng Ni,		
745	and Ting Liu. 2021. Co-gat: A co-interactive graph		
746	attention network for joint dialog act recognition		
747	and sentiment classification. In <i>Proceedings of the</i>		
748	<i>AAAI conference on artificial intelligence</i> , volume 35,		
749	pages 13709–13717.		
750	Kaisong Song, Lidong Bing, Wei Gao, Jun Lin, Lujun		
751	Zhao, Jiancheng Wang, Changlong Sun, Xiaozhong		
752	Liu, and Qiong Zhang. 2019. Using customer ser-		
753	vice dialogues for satisfaction analysis with context-		
754	assisted multiple instance learning. In <i>Proceedings of</i>		
755	<i>the 2019 conference on empirical methods in natural</i>		
756	<i>language processing and the 9th international joint</i>		
757	<i>conference on natural language processing (EMNLP-</i>		
758	<i>IJCNLP)</i> , pages 198–207.		