

Uni-ETOD: User-Need-Driven Chain of Thought Framework for Fully End-to-end Task-oriented Dialogue System

Anonymous ACL submission

Abstract

Fully End-to-End Task-Oriented Dialogue Systems (Fully ETOD) retrieve knowledge from a knowledge base in a differentiable manner and generate responses using a language model generator without the need for modular training. However, Fully ETOD faces some challenges. During the retrieval process, the retriever retrieves the knowledge base in a black-box manner, making it difficult for the generator to differentiate the large amount of knowledge obtained by the retriever. This leads to a degradation in the quality of the responses and the trustworthiness of the system. Moreover, as the size of the knowledge base grows, it may exacerbate the risk of this problem. To address this challenge, we first design a dataset for Fully ETOD based on large-scale knowledge bases called FakeRest to solve the scarcity of annotated dialogue data based on large-scale knowledge bases. We also propose a User-need-driven Chain of Thought Framework (Uni-ETOD) for Fully ETOD, which aims to guide LLMs to gradually understand users' thought processes and improve the quality of responses in Fully ETOD. We use ChatGPT, Gemini, Llama3, Mistral, and ChatGLM as the backbone models of the system. On FakeRest, we comprehensively evaluate the capability of each step of Uni-ETOD. The results show that Uni-ETOD will help LLMs better distinguish the retrieved knowledge and enhance the credibility and interpretability of the whole system.

1 Introduction

Task-oriented dialogue (TOD) systems can accomplish specific tasks, such as booking restaurant reservations and providing transportation navigation, through user interaction and leveraging an external knowledge base. Traditional TOD systems follow a pipeline approach and consist of four interconnected modular components (Qin et al., 2020; Jacqmin et al., 2022; Hosseini-Asl et al., 2020).



Figure 1: A sample demonstration of Fully ETOD. In the first round, the system did not select all the entities that met the needs. In the second round, the system selected irrelevant attributes for the response. We mark the correct entities and attributes in green. Ignored or incorrect entities and attributes are labeled in red.

The Modularly End-to-End Task-oriented Dialogue System (Modularly EToD) trains all components in an end-to-end manner while optimizing their parameters. In contrast to traditional TOD and Modularly EToD, the Fully Task-oriented End-to-End Dialogue System (Fully ETOD) (Eric and Manning, 2017) encodes knowledge bases (KBs) and uses neural networks to query the KBs in a differentiable manner. Fully ETOD generates a system response directly given only the dialogue history and the corresponding knowledge base. Therefore, it has received great attention from both academia and industry. Although Fully ETOD exhibits excellent data scalability, efficiently retrieving the desired entities and attributes poses a challenge due to the abundance of irrelevant knowledge in

058 the retrieved results. This issue hampers the genera- 110
059 tor’s ability to distinguish between retrieved entities 111
060 effectively and extract valid information. 112

061 Existing Fully ETODs tend to follow the 113
062 retrieval-generation paradigm. Retrieval- 114
063 augmented generation (RAG) (Lewis et al., 2020; 115
064 Ren et al., 2021; Singh et al., 2021) improves the 116
065 quality and relevance of the generated text and 117
066 achieves positive results in knowledge-intensive 118
067 tasks. Q-TOD (Tian et al., 2022) efficiently 119
068 implements retrieval-enhanced generation on Fully 120
069 ETOD, thus alleviating the domain adaptation 121
070 problem. MK-TOD (Shen et al., 2023) raises the 122
071 problem of mismatch between the retrieval and 123
072 generation processes, and incorporates a variety of 124
073 meta-knowledge to guide the generator to increase 125
074 the retrieval utilization of the results. However, for 126
075 the generator, the retrieval process of the retriever 127
076 remains a black-box process. This means that 128
077 although a retriever can provide retrieval results 129
078 with high recall, the results still contain a large 130
079 number of irrelevant entities and attributes. It is 131
080 difficult for the generator to differentiate between 132
081 these retrieved entities and to select the attributes 133
082 of the entities that meet the needs, as shown in 134
083 Figure 1. This is referred to as the low precision 135
084 problem of the retrieval process. Meanwhile, the 136
085 black-box and low-precision retrieval process is 137
086 difficult to analyze, which damages the user’s trust 138
087 (Qin et al., 2023). We consider the low precision 139
088 problem and non-interpretability of the retrieval 140
089 process as bottlenecks in the existing Fully ETOD. 141
090 Another core challenge with existing Fully ETOD 142
091 is the lack of annotated dialogue data based on 143
092 large-scale knowledge bases (Qin et al., 2023). 144
093 Existing Fully ETOD datasets are often based on 145
094 small knowledge bases or modified from existing 146
095 datasets. Due to not being tailored for large-scale 147
096 datasets, existing Fully ETOD datasets frequently 148
097 encounter discrepancies between responses and 149
098 annotated knowledge. 150

099 Recent developments in large language model- 151
100 ing provide us with solutions to the above problems. 152
101 Our paper introduces a User-need-driven Chain of 153
102 Thought Framework for Fully ETOD (Uni-ETOD). 154
103 Uni-ETOD aims to improve the precision and inter- 155
104 pretability of the retrieval process, and thus the 156
105 quality of the responses. The framework consists 157
106 of three steps: 1) Retrieval Based on User Needs: 158
107 Based on the dialogue context, LLMs will generate 159
108 the user’s needs and use the embedding model to re-
109 trieve the most relevant knowledge in a large-scale

knowledge base. 2) Knowledge Refinement: Based 110
on the retrieval results obtained by the retriever, 111
the LLMs will filter the entities and attributes that 112
match the user’s needs, resulting in a more accu- 113
rate retrieval result. The results can be directly 114
displayed to the user as part of the response. 3) 115
Response Based on Refined Knowledge: Based 116
on the retrieval results with higher precision, the 117
generator will make a more credible response. 118

119 In addition, we propose an automated method 120
for constructing dialog data based on a large-scale 121
knowledge base. We utilize LLMs to simulate 122
restaurant scenarios and construct FakeRest. Fak- 123
eRest is specifically tailored for Fully ETOD of 124
large-scale knowledge bases, containing detailed 125
annotated data to enhance the knowledge base re- 126
trieval capability of LLMs. Refer to chapter 3 for 127
more details.

128 We apply Uni-ETOD to several LLMs, includ- 129
ing two closed-source models, ChatGPT (Brown 130
et al., 2020), Gemini (Team et al., 2023), and three 131
open-source models, Llama3 (AI@Meta, 2024), 132
Mistral (Jiang et al., 2023), and ChatGLM (Du 133
et al., 2022). On the FakeRest dataset, we utilize 134
the task to evaluate the enhancement of Uni-ETOD 135
on the retrieval process and the response process, 136
respectively. The experimental results show that 137
Uni-ETOD can effectively alleviate the low preci- 138
sion problem in the retrieval process, and enhance 139
the quality of responses and interpretability.

2 Related Work 140

2.1 Fully End-to-End Task-oriented Dialog 141

142 We use whether or not we query KBs as APIs us- 143
ing beliefs as a criterion to differentiate between 144
modular ETODs and Fully ETODs. Fully ETODs 145
retrieve knowledge bases in a differentiable way. 146
We classify existing Fully ETODs as being cate- 147
gorized into two types. First, the knowledge base 148
is stored in the model parameters, and the system 149
retrieves it implicitly and generates a response to 150
the user. GPT-KE (Madotto et al., 2020) learns 151
knowledge base embedding through data augmen- 152
tation and responding to users. ECO (Huang et al., 153
2022) guides the response generation in the end-to- 154
end system by autoregressively generating entities. 155
However, such an approach mixes the retrieval and 156
response processes and inevitably generates low- 157
confidence generation results, especially when the 158
size of the knowledge base becomes large.

159 Second, the system explicitly retrieves KBs

through a retriever. DialoKG (Rony et al., 2022) selects relevant triples by graph embedding. Q-TOD (Tian et al., 2022) utilizes a rewritten query in combination with the RAG technique to improve retrieval performance. MAKER (Wan et al., 2023) queries entities and attributes separately through two retrievers. MK-TOD (Shen et al., 2023) mitigates the misalignment between the retriever and the generator by combining meta-knowledge. Although retrieving the KB explicitly provides the retrieved results compared to the first approach, the black-box retrieval process still limits the interpretability. Moreover, such a retrieval process will inevitably introduce a large number of irrelevant entities and attributes. Our work aims to utilize the power of LLMs to alleviate the interpretability and low precision problems of the retrieval process.

2.2 Large Language Models for ETOD

Recently, LLMs have achieved great success and demonstrated extraordinary text generation and reasoning capabilities (Suzgun et al., 2023; Pu and Demberg, 2023; Kojima et al., 2022). Unlike small language models that have difficulty solving complex problems, LLMs can solve complex problems with various prompting strategies (Zhao et al., 2023). This is based on the amazing ability that LLMs show in multi-hop reasoning. Wei et al. (Wei et al., 2022) investigated the Chain of Thought (CoT) prompting technique in LLMs by inducing the model to generate intermediate steps to improve the precision of answers. Meanwhile, based on the powerful contextual learning capability of LLMs, many existing works combine LLMs with traditional TOD systems. They generally use the zero-shot or few-shot approach to explore the capability of LLMs applied to individual modules (Pan et al., 2023; Heck et al., 2023; Hudeček and Dušek, 2023; Parikh et al., 2023). However, there is a gap in the work on applying LLMs to Fully ETOD. The lack of a Fully ETOD dataset and training paradigm based on a large-scale knowledge base is a hindrance.

3 FakeRest: A Fully ETOD Dataset for the Large-Scale Database

3.1 Why build the FakeRest dataset?

To address the challenge of scarce labeled data faced by Fully ETOD, we propose a construction method and construct FakeRest, a dataset of scheduled restaurant scenarios designed for Fully ETOD

based on a large-scale knowledge base. We utilize LLMs to simulate the dialog scenarios between the user and the system in the restaurant. Based on a predefined user need path, the user LLM and the system LLM will have multiple rounds of dialog until the system finds (or fails to find) the only restaurant that matches the user’s need. In this process, we will record the thought and retrieval process of the user and the system in detail as annotated data.

FakeRest has the following advantages:

1) **Based on Large-Scale Knowledge Base:** Unlike the existing Fully ETOD datasets, which are based on a small-scale knowledge base. FakeRest’s knowledge base contains entities for 120 different restaurants. Similar to the format of CamRest (Rojas-Barahona et al., 2016), each restaurant in the knowledge base contains seven attribute values. Four are private attributes (name, phone, address, postcode) for each restaurant, and the values of the private attributes are completely different. Including three public attributes (area, food, price range), the public attributes can match the user’s needs.

2) **Data Consistency:** The existing Fully ETOD datasets annotate the corresponding knowledge base based on the dialogue content of the original dataset. This could lead to inconsistencies between the dialogue content and the annotated knowledge base, or the system may fail to respond based on all the entities that match the user’s needs. In contrast, during each dialogue round in FakeRest, the service system will retrieve all the restaurants that match the user’s needs from the knowledge base of 120 restaurants and respond to the user. Such an approach ensures the consistency of the annotated knowledge and responses.

3) **More Detailed Annotation:** In addition annotating the user’s and system’s utterances, we provided detailed annotations of the user’s and system’s thought processes, including the user’s needs, all entity IDs, and attribute values that aligned with the user’s needs for the round. To maintain data diversity and balance, we established various need paths and target restaurants for each dialogue. We also created scenarios where no restaurants in the knowledge base matched the user’s needs.

3.2 How to build the FakeRest dataset?

We propose a method to automatically construct dialog data for large-scale knowledge bases and make an attempt with a subscription restaurant scenario. First, we construct a knowledge base of 120 restau-

rants using ChatGPT (Brown et al., 2020). Then, we designed different need paths based on different attributes of the restaurants in the knowledge base. Each need path corresponds to one LLM user. We design 900 users, among which 720 users found the needed restaurants and 180 users did not. Finally, based on the need paths, we use templates to construct query statements to find all the entities and attributes that match the user’s needs, and as part of the annotations. The LLM user and the LLM system will generate multiple rounds of dialogs based on these annotations. We describe these three steps in detail next. Please refer to the Appendix A for prompts.

Constructing the Knowledge Base We defined 18 public attributes (5 area, 10 food, 3 price range) and simulated 120 restaurants based on different combinations of public attributes. It is worth noting that we could have simulated 150 restaurants with slightly different public attributes. However, in order to generate scenarios that did not align with the user’s needs, we randomly removed 30 restaurants. Subsequently, we used ChatGPT to create unique private attributes for each restaurant, which were then reviewed and refined manually. As a result, we obtained 120 distinct restaurants.

Constructing the Need Path Similar to the process of constructing restaurants, we construct the user’s need path by combining various public attributes. We use food as the starting point of the path and choose either area or price range as the subsequent step. For instance, if the user is looking for a Korean restaurant in the eastern part of the city, the need path for this round is [Korean(food), east(area)]. Each path will end in the final round by finding the only restaurant that meets the need or by not finding any suitable restaurant. Once the restaurant is identified, the need path will take the private attribute as the next step (e.g., postcode, phone). We assign the corresponding LLM users based on the 900 need paths and utilize the interactions between each user and the system as training data.

Constructing Detailed Annotated Data We transform the need path into a fluent sentence as a user need and allow the LLM user to express the current need based on the dialog context and the user need. For instance, "You are looking for a Korean restaurant in the east area.". Unlike the subsequent methods that utilize RAGs for retrieval, we convert the need path into a deterministic search statement to search a list of all restaurants in the

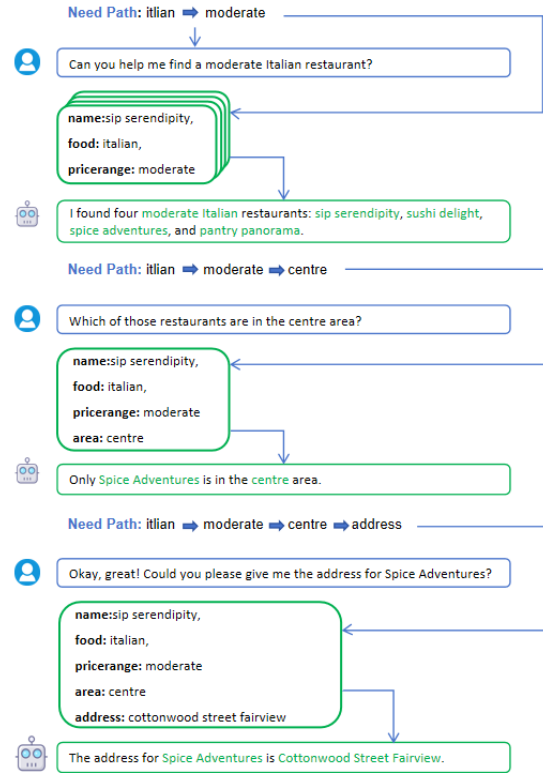


Figure 2: Figure of the process of constructing FakeRest. Based on the set user’s need path, the user LLM and system LLM simulate the restaurant reservation scenario.

knowledge base that meet the need. For more precise labeling, we label the IDs of the retrieved entities and the values of the attributes of the current round of needs, rather than all the entity information. For example, {"name": "taste discoveries", "food": "korean", "area": "east"}. The system will respond to the LLM user with detailed labeled information based on the context of the dialogue.

3.3 Dataset Statistics

In Table 1 we summarize the statistics for FakeRest. We use the partitioning of training/test, where #Entities denotes the number of entities in the knowledge base and #Attributes denotes the number of attribute values in the knowledge base.

Dataset	#Dialogues	#Utterances	#Entites	#Attributes
FakeRest	765/135	4252/752	120	498

Table 1: Dataset statistics of the FakeRest dataset.

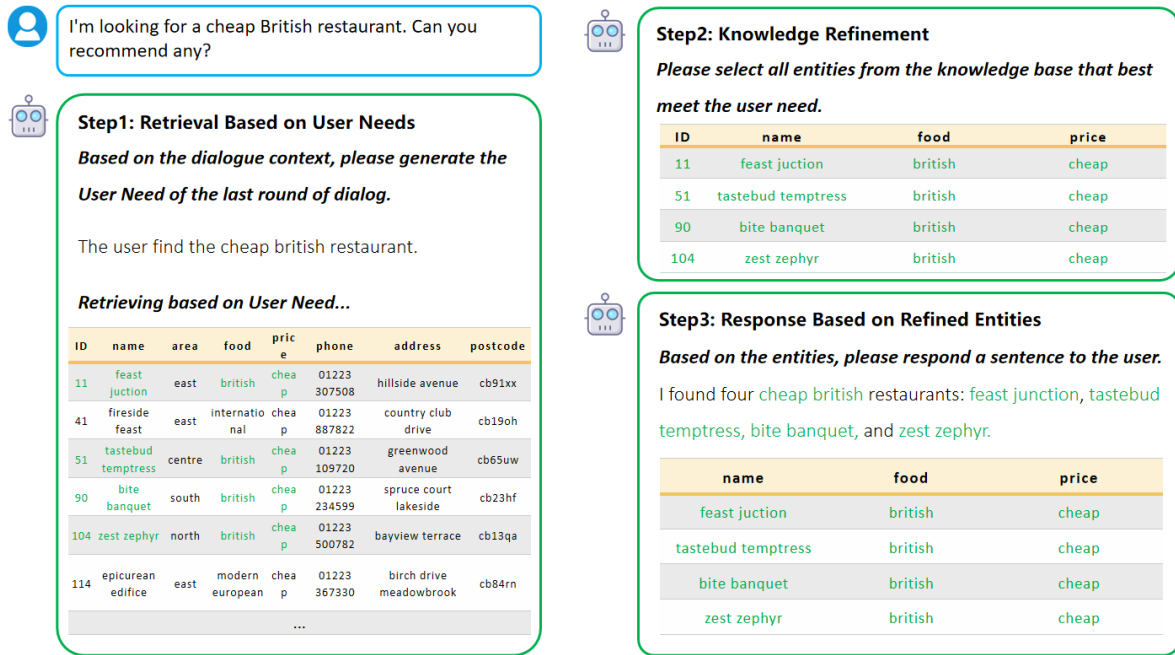


Figure 3: An illustration of a Uni-ETOD. In each round of dialogue, Uni-ETOD retrieves and refines knowledge according to the user’s needs. To better represent the knowledge base information, we use a graph instead of a JSON list. Finally, Uni-ETOD will be able to return structured knowledge that is interpretable after refinement (e.g., a list of restaurants that meet the user’s needs) and a paragraph of highly reliable response.

4 Uni-ETOD: User-Need-Driven Chain of Thought Framework for Fully ETOD

The Fully ETOD task can be defined as given a dialog history H and an associated knowledge base KB . $H = (u_1, s_1), (u_2, s_2), \dots, (u_{n-1}, s_{n-1}), u_n$, where u_n and s_n denote the n -th round of user utterances and system utterances, respectively. The $KB = (e_1, e_2, \dots, e_m)$, $e_m = (a_1, a_2, \dots, a_k)$, where e_m and a_k denote entities and attributes in the knowledge base. The purpose of the system agent is to predict the system response s_n , denoted as S .

In recent years, LLMs have changed the paradigm of natural language, showing strong multi-hop reasoning capabilities (Zhao et al., 2023). Inspired by CoT (Wei et al., 2022), we consider replacing the traditional retrieval-generation paradigm with a multi-step reasoning framework. We aim to assist LLMs to understand the user’s needs step by step and explain an understandable retrieval process based on the user’s needs. The retrieval results will show which entities and attribute values match the user’s needs and ultimately provide a high-quality response. We will detail the three steps of Uni-ETOD and the required prompt templates.

4.1 Retrieval Based on User Needs

In this step, we first construct a dataset corresponding to user needs and entities in the knowledge base using the training set of FakeRest. Since FakeRest has detailed annotations, we construct a sentence pair dataset of user needs and entities for fine-tuning the embedding model. The dataset can be represented as $D = (un_1, [e_1, e_2, \dots, e_i]), \dots, (un_i, [e_1, e_2, \dots, e_j])$, where un and e represent the user’s needs and the entities in the knowledge base, respectively. We train our embedding model using all positive samples to compute similarity scores and cross-entropy loss of labels. See the Appendix B for experimental details.

Then we use LLMs to summarize the user need based on the dialogue history H . We use the following template:

Based on the 'Dialog Context', please generate the 'User Need' of the 'Last Round' of dialog.
 Dialog Context: [Dialogue History H]
 Last Round: [u_n]

Given the dialogue history H , we allow the LLMs to summarize the user’s needs for the last round of dialogue u_n . This step can be represented as Equation 1, where U is a concise sentence representing the user’s needs.

$$U = \operatorname{argmax} P(y|H, u_n) \quad (1)$$

Finally, we use the U to retrieve the top-k entities in the knowledge base that match the user’s needs, which can be represented as:

$$E_{top-k} = P(U, KB) \quad (2)$$

4.2 Knowledge Refinement

In this step, we allow the LLMs to refine the retrieved knowledge based on the user’s needs. The LLMs filter out entities and attributes that align with the user’s needs from the retrieved knowledge. This step is designed to address the low precision problem of the retrieved results, aiming to enhance the credibility of the retrieved knowledge. Furthermore, we allow LLMs to provide additional explanations for the process and present high-precision knowledge to the user, improving the system’s interpretability. We follow the template below:

Please select all ‘Entities’ from the ‘Knowledge Base’ that best meet the user needs.
 Dialogue Context: [Dialogue History H]
 Knowledge Base: [Retrieved Entities E]

Given the dialog history H and the retrieval result E from the previous step, we guide the LLMs to select the knowledge that meets the user’s needs. The refined knowledge contains only the entities and attributes that match the user’s needs. This step can be represented as Equation 3, where $E_{refined}$ is the set of refined entities.

$$E_{refined} = \operatorname{argmax} P(y|H, E) \quad (3)$$

4.3 Response Based on Refined Knowledge

In this step, LLMs provide high quality responses to customers based on retrieval results with high precision and recall. We use the following template:

Based on the ‘Entities’, please respond a sentence to the user.

Dialog Context: [Dialog History H]

Entities: [Refined Entities $E_{refined}$]

Given the dialog history H and the refined retrieval entities $E_{refined}$, we guide the LLMs to generate the final response. This step can be represented as Equation 4. S represents the final system response.

$$S = \operatorname{argmax} P(y|H, E_{refined}) \quad (4)$$

5 Experiments

We first explain the problem of low precision in the retrieval process. To comprehensively evaluate the effectiveness of Uni-ETOD, we will assess the retrieval and response processes of Fully ETOD separately.

5.1 Low Precision of the Retrieval Process

As the scale of the knowledge base increases, it becomes increasingly challenging for the retriever to find all the entities that match the user’s needs from a large number of entities. Many works enhance the recall of the retrievers through more parameterized encoders and sophisticated training methods. However, the retrieval results inevitably contain numerous irrelevant entities and attributes. Since the entities in the knowledge base are often very similar and there are usually very few entities that match the user’s needs, this results in low precision problems in the retrieval process. These irrelevant entities and attributes will be directly utilized as the basis for the generator to make a response, which will destroy the credibility of the response.

Shen et al. (Shen et al., 2023) proposed the retrieval-generation misalignment problem, which refers to the inconsistency between the recall rate of the retrieval process and the response quality. We believe that the problem of low precision in retrieval results is the main reason for this phenomenon. This is because retrieval results that contain a large number of irrelevant entities will be more dependent on the generator’s capabilities. Additionally, as illustrated in Table 2, even a high recall retriever exacerbates the low precision problem when increasing the k value to handle larger retrievals, impacting the response precision. The low precision problem may hinder the development of Fully ETOD based on large-scale knowledge bases,

as more efficient retrievers may not effectively address this issue. Our approach addresses this problem by leveraging the reasoning abilities of LLMs through stepwise reasoning. In our results, Uni-ETOD notably enhances the precision of retrieval results and demonstrates improved alignment between the retrieval and generation processes.

k	Precision@10	Recall@10	Entity Precision	Entity Recall
5	27.5	80	63.9	70.3
10	16.1	94	63.7	70.5
15	11.1	97.2	62.1	73.4
20	8.5	99.1	60.1	73.3

Table 2: Trends in precision and recall during retrieval and response as k increases

5.2 Overall Results on Retrieval Process

We used bge-large-en-v1.5 (Xiao et al., 2023) as the base model and also compared the retrieval results of an embedded model that utilizes fine-tuning of user needs (without Knowledge Refinement) as well as UniETOD. To evaluate the retrieval capability of Fully ETOD, the recall of the retrieval process is evaluated in addition to Recall@k (Shen et al., 2023). We also propose Precision@k and F1@k to evaluate the precision and comprehensive performance of the retrieval process.

Model	Precision@10	Recall@10	F1@10
BGE-Large	10.4	60.4	17.7
BGE-Large+User Need	16.9	98.9	29.0
Uni-ETOD(ChatGLM)	79.4	85.9	82.5
Uni-ETOD(Llama3)	96.0	96.7	96.4
Uni-ETOD(Mistral)	97.2	94.1	95.7
Uni-ETOD(Gemini)	63.2	92.9	75.3
Uni-ETOD(ChatGPT)	70.0	89.0	78.3

Table 3: Experimental results of the retrieval process on FakeRest dataset.

As shown in Table 3, we observe that the fine-tuned embedding model, based on FakeRest’s detailed annotation, performs well in achieving very high recall results for short user needs. The fine-tuned embedding model enhances recall by 38.5 percent under the top-10 retrieval setting. However, the retrieval process is hindered by the low precision problem due to the constraints of the large-scale knowledge base. This problem is well mitigated by Uni-ETOD, which improves precision at the expense of a minimal reduction in recall. For instance, Uni-ETOD (Llama3) decreases the recall by 2.2 percent but improves the precision by 79.1 percent. Although ChatGPT and Gemini were unable to fine-tune the Knowledge Refined step, both

models still improved the precision of the retrieval results in the zero-shot setting. The results demonstrate that Uni-ETOD improves the precision of the retrieval process, thereby elevating the overall quality of Fully ETOD. Moreover, the user-need-based retrieval process in Uni-ETOD notably enhances the recall of the base model, achieving 98.9 percent at the top-10 setting.

5.3 Overall Results on Response Process

In this section, we show all the implementation details and all the experimental results of the LLMs in the response process.

5.3.1 Implementation Details

We use 5 LLMs for our experiments, including 2 closed-source models and 3 open-source models. Specifically, we use GPT3.5 (gpt-3.5-turbo) (Brown et al., 2020) from the OpenAI API and Gemini (gemini-1.5-flash) (Team et al., 2023)] from Google Gemini API. Additionally, we include three open-source models in our experiments: ChatGLM3 (chatglm3-6B) (Du et al., 2022) and Llama3 (meta-llama3-8B-instruct) (AI@Meta, 2024) and Mistral (mistral-7B-instruct-v0.2) (Jiang et al., 2023). The temperature is set to 0.1, while all other hyperparameters are set to default values.

We use LoRA (Hu et al., 2021) to fine-tune all LLMs. The LLMs are fine-tuned on a 24G NVIDIA 3090. We set training batch size to 1, epoch number to 3, learning rate to 5e-5, and warm-up steps to 20.

We compare the improvement of Uni-ETOD on the overall performance and reliability of LLMs in the response process by using retrieve-generate paradigm (base) and retrieve+zero-shot CoT (zeroCoT) (Kojima et al., 2022) as a baseline. To evaluate the quality of Fully ETOD’s responses, we use the BLEU (Papineni et al., 2002), Entity F1 (Eric and Manning, 2017) metrics to assess the consistency of responses and the generator’s ability to respond with correct knowledge.

5.3.2 Results on Zero-shot Reasoning and Supervised Fine-tuning

Since Gemini and ChatGPT are not fine-tunable, we show performance with zero shots as well as results with the fine-tuned Llama3 as the knowledge refinement component(+KR). On Llama3, ChatGLM, and Mistral we show performance with zero-shot inference and fine-tuning settings on the response process, respectively.

Model	BLEU	Entity Precision	Entity Recall	Entity F1	
ChatGPT	base	24.1	67.5	69.9	68.7
	zeroCoT	19.0	57.8	85.3	68.9
	Uni-ETOD	26.2	70.9	69.9	70.4
	Uni-ETOD+KR(finetuned)	33.3	84.8	88.3	86.5
	base	37.3	73.3	81.0	76.9
Gemini	zeroCoT	38.9	75.9	83.0	79.3
	Uni-ETOD	43.0	77.3	89.3	82.9
	Uni-ETOD+KR(finetuned)	55.6	94.3	94.8	94.6
Llama3	base	10.2	63.7	70.5	66.9
	zeroCoT	19.1	65.4	77.7	71.0
	Uni-ETOD	15.2	76.0	82.4	79.0
	Uni-ETOD(finetuned)	56.0	94.5	94.1	94.3
	base	6.3	43.1	76.2	55.1
Mistral	zeroCoT	6.5	40.2	75.7	52.5
	Uni-ETOD	12.4	69.1	88.7	77.7
	Uni-ETOD(finetuned)	55.0	95.5	90.3	92.8
ChatGLM	base	8.1	39.0	73.7	51.0
	zeroCoT	6.7	36.6	75.2	49.2
	Uni-ETOD	21.8	76.2	87.7	81.6
	Uni-ETOD(finetuned)	43.9	91.8	86.6	89.1

Table 4: Experimental results of the response process on FakeRest dataset.

As shown in Table 4, Uni-ETOD significantly improves the quality of responses compared to the baseline method. Uni-ETOD effectively mitigates the low precision problem of the retrieval process, which is reflected in the quality of responses. Uni-ETOD can better utilize the retrieval results with high recall and precision to generate more comprehensive and high confidence responses. In particular, the fine-tuned Uni-ETOD (Llama3) achieves a BLEU as high as 56 percent and Entity F1 as high as 94.3, outperforming ChatGPT and Gemini with zero-shot setting.

The experiments demonstrate that Uni-ETOD consistently enhance the performance of LLMs in the response process of Fully ETOD. We argue that Uni-ETOD can effectively stimulate LLMs’ multi-hop reasoning in Fully ETOD. By guiding LLMs to gradually understand users’ needs, Uni-ETOD can provide users with higher-quality and more credible responses.

5.4 Ablation Study

First, we evaluate the role of each step in Uni-ETOD in the retrieval process. Due to computational resource constraints, we performed ablation experiments on the fine-tuned Llama3 on FakeRest.

Model	Precision@10	Recall@10	F1@10
Uni-ETOD	96.0	96.7	96.4
w/o KR	17.0	98.9	29.0
w/o UR	97.8	91.5	94.6
w/o KR & UR	10.4	60.4	17.7

Table 5: Ablation study of the retrieval process on FakeRest dataset.

The table 5 demonstrates the impact of Retrieval Based on User Needs (RU), and Knowledge Refinement (KR) on the Fully ETOD retrieval process.

Our RU step contains fine-tuning of the embedding model. For a fair comparison, we fine-tuned dialogue history as the query for “w/o RU” and “w/o RU & UR”, and Uni-ETOD still achieved better results. The results show that the RU step can effectively improve the recall of the original paradigm retrieval process. Additionally, KR can effectively alleviate the low precision problem in the retrieval process.

We evaluate the role of each step in Uni-ETOD in the response process. We perform ablation experiments on fine-tuned Llama3, ChatGLM, and Mistral on FakeRest.

Model	BLEU	Entity Precision	Entity Recall	Entity F1	
Llama3	Uni-ETOD	56.0	94.5	94.1	94.3
	w/o KR	40.7	85.0	88.9	85.0
	w/o UR	51.5	94.2	90.7	92.4
	w/o KR & UR	34.9	82.6	80.4	81.5
Mistral	Uni-ETOD	55.0	95.5	90.3	92.8
	w/o KR	43.9	85.6	88.9	87.2
	w/o UR	52.0	95.1	88.5	91.7
	w/o KR & UR	36.4	85.5	79.4	82.3
ChatGLM	Uni-ETOD	43.9	91.8	86.6	89.1
	w/o KR	29.8	55.2	83.8	66.5
	w/o UR	39.7	91.5	82.8	86.9
w/o KR & UR	27.0	57.4	76.0	65.4	

Table 6: Ablation study of the response process on FakeRest dataset.

The table 6 demonstrates the effects of Retrieval Based on User Needs (RU) and Knowledge Refinement (KR) on the Fully ETOD response process. The enhancement brought in the retrieval process is also shown in the response process. The results show that both RU and KR can effectively mitigate the low precision problem, improving the credibility, and overall quality of the responses.

6 Conclusion

In this paper, we aim to address the problems of low precision and poor interpretability in Fully ETOD. We propose a user-need-driven CoT framework (Uni-ETOD), which allows LLMs to gradually understand user needs and generate high-quality responses through multi-step reasoning. Experimental results demonstrate that Uni-ETOD effectively alleviates the low precision problem and offers users a more explanatory retrieval process and more reliable responses. Furthermore, we present a technique for automatically generating dialog data based on large-scale knowledge bases and constructing FakeRest, a dialogue dataset for restaurant scenarios.

596 Limitations

597 There are two limitations of this paper that de-
598 serve a deeper examination. First, we have not ex-
599 plored the fine-tuning methods and sampling tech-
600 niques for embedding models in depth. Second,
601 the method in this paper can be fine-tuned to adapt
602 to various scenarios by automatically generating
603 dialog data. However, our method is still not an
604 autonomous learning method to adapt the system
605 to new scenarios through user interaction.

606 References

607 AI@Meta. 2024. [Llama 3 model card](#).

608 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
609 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
610 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
611 Askell, et al. 2020. Language models are few-shot
612 learners. *Advances in neural information processing*
613 *systems*, 33:1877–1901.

614 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding,
615 Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm:
616 General language model pretraining with autoregres-
617 sive blank infilling. In *Proceedings of the 60th An-
618 nual Meeting of the Association for Computational*
619 *Linguistics (Volume 1: Long Papers)*, pages 320–335.

620 Mihail Eric and Christopher D Manning. 2017. Key-
621 value retrieval networks for task-oriented dialogue.
622 *arXiv preprint arXiv:1705.05414*.

623 Michael Heck, Nurul Lubis, Benjamin Matthias Ruppik,
624 Renato Vukovic, Shutong Feng, Christian Geishausser,
625 Hsien chin Lin, Carel van Niekerk, and Milica
626 Gavsic. 2023. [Chatgpt for zero-shot dialogue state](#)
627 [tracking: A solution or an opportunity?](#) *ArXiv*,
628 [abs/2306.01386](#).

629 Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu,
630 Semih Yavuz, and Richard Socher. 2020. A simple
631 language model for task-oriented dialogue. *Advances*
632 *in Neural Information Processing Systems*, 33:20179–
633 20191.

634 J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan
635 Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu
636 Chen. 2021. [Lora: Low-rank adaptation of large](#)
637 [language models](#). *ArXiv*, [abs/2106.09685](#).

638 Guanhuan Huang, Xiaojun Quan, and Qifan Wang.
639 2022. [Autoregressive entity generation for end-to-](#)
640 [end task-oriented dialog](#). In *Proceedings of the 29th*
641 *International Conference on Computational Linguis-*
642 *tics*, pages 323–332, Gyeongju, Republic of Korea.
643 International Committee on Computational Linguis-
644 tics.

645 Vojtěch Hudeček and Ondřej Dušek. 2023. Are large
646 language models all you need for task-oriented dia-
647 logue? In *Proceedings of the 24th Annual Meeting*

of the Special Interest Group on Discourse and Dia-
648 *logue*, pages 216–228. 649

Léo Jacqmin, Lina M. Rojas Barahona, and Benoit
650 Favre. 2022. [“do you follow me?”: A survey of](#)
651 [recent approaches in dialogue state tracking](#). In *Pro-*
652 *ceedings of the 23rd Annual Meeting of the Special*
653 *Interest Group on Discourse and Dialogue*, pages
654 336–350, Edinburgh, UK. Association for Computa-
655 tional Linguistics. 656

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-
657 sch, Chris Bamford, Devendra Singh Chaplot, Diego
658 de las Casas, Florian Bressand, Gianna Lengyel, Guil-
659 laume Lample, Lucile Saulnier, et al. 2023. Mistral
660 7b. *arXiv preprint arXiv:2310.06825*. 661

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-
662 taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-
663 guage models are zero-shot reasoners. *Advances in*
664 *neural information processing systems*, 35:22199–
665 22213. 666

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
667 Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
668 rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-
669 täschel, et al. 2020. Retrieval-augmented generation
670 for knowledge-intensive nlp tasks. *Advances in Neu-*
671 *ral Information Processing Systems*, 33:9459–9474. 672

Ilya Loshchilov and Frank Hutter. 2017. [Decoupled](#)
673 [weight decay regularization](#). In *International Confer-*
674 *ence on Learning Representations*. 675

Andrea Madotto, Samuel Cahyawijaya, Genta Indra
676 Winata, Yan Xu, Zihan Liu, Zhaojiang Lin, and Pas-
677 cale Fung. 2020. [Learning knowledge bases with](#)
678 [parameters for task-oriented dialogue systems](#). In
679 *Findings of the Association for Computational Lin-*
680 *guistics: EMNLP 2020*, pages 2372–2394, Online.
681 Association for Computational Linguistics. 682

Wenbo Pan, Qiguang Chen, Xiao Xu, Wanxiang Che,
683 and Libo Qin. 2023. A preliminary evaluation of
684 chatgpt for zero-shot dialogue understanding. *arXiv*
685 *preprint arXiv:2304.04256*. 686

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-
687 Jing Zhu. 2002. [Bleu: a method for automatic evalu-](#)
688 [ation of machine translation](#). In *Annual Meeting of*
689 *the Association for Computational Linguistics*. 690

Soham Parikh, Quaizar Vohra, Prashil Tumbade, and
691 Mitul Tiwari. 2023. Exploring zero and few-shot
692 techniques for intent classification. *arXiv preprint*
693 *arXiv:2305.07157*. 694

Dongqi Pu and Vera Demberg. 2023. [Chatgpt vs human-](#)
695 [authored text: Insights into controllable text summa-](#)
696 [rization and sentence style transfer](#). In *Annual Meet-*
697 *ing of the Association for Computational Linguistics*. 698

Libo Qin, Wenbo Pan, Qiguang Chen, Lizi Liao, Zhou
699 Yu, Yue Zhang, Wanxiang Che, and Min Li. 2023.
700 [End-to-end task-oriented dialogue: A survey of tasks,](#)
701 [methods, and future directions](#). In *Proceedings of the*
702

703				
704				
705				
706	Libo Qin, Xiao Xu, Wanxiang Che, and Ting Liu. 2020.			
707	AGIF: An adaptive graph-interactive framework for			
708	joint multiple intent detection and slot filling.			
709	In <i>Findings of the Association for Computational Lin-</i>			
710	<i>guistics: EMNLP 2020</i> , pages 1807–1816, Online.			
711	Association for Computational Linguistics.			
712	Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao,			
713	QiaoQiao She, Hua Wu, Haifeng Wang, and Ji-Rong			
714	Wen. 2021. RocketQAv2: A joint training method			
715	for dense passage retrieval and passage re-ranking.			
716	In <i>Proceedings of the 2021 Conference on Empiri-</i>			
717	<i>cal Methods in Natural Language Processing</i> , pages			
718	2825–2835, Online and Punta Cana, Dominican Re-			
719	public. Association for Computational Linguistics.			
720	Lina Maria Rojas-Barahona, Milica Gaić, Nikola Mrk-			
721	sić, Pei hao Su, Stefan Ultes, Tsung-Hsien Wen,			
722	Steve J. Young, and David Vandyke. 2016. A			
723	network-based end-to-end trainable task-oriented di-			
724	alogue system.			
725	In <i>Conference of the European Chapter of the Association for Computational Linguistics</i> .			
726	Md. Rony, Ricardo Usbeck, and Jens Lehmann. 2022.			
727	Dialokg: Knowledge-structure aware task-oriented			
728	dialogue generation.			
729	<i>ArXiv</i> , abs/2204.09149.			
730	Weizhou Shen, Yingqi Gao, Canbin Huang, Fanqi			
731	Wan, Xiaojun Quan, and Wei Bi. 2023. Retrieval-			
732	generation alignment for end-to-end task-oriented			
733	dialogue system.			
734	In <i>Proceedings of the 2023 Confer-</i>			
735	<i>ence on Empirical Methods in Natural Language</i>			
736	<i>Processing</i> , pages 8261–8275, Singapore. Associa-			
737	tion for Computational Linguistics.			
738	Devendra Singh, Siva Reddy, Will Hamilton, Chris			
739	Dyer, and Dani Yogatama. 2021. End-to-end train-			
740	ing of multi-document reader and retriever for open-			
741	domain question answering. <i>Advances in Neural</i>			
742	<i>Information Processing Systems</i> , 34:25968–25981.			
743	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Se-			
744	bastian Gehrmann, Yi Tay, Hyung Won Chung,			
745	Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny			
746	Zhou, and Jason Wei. 2023. Challenging BIG-bench			
747	tasks and whether chain-of-thought can solve them.			
748	In <i>Findings of the Association for Computational Lin-</i>			
749	<i>guistics: ACL 2023</i> , pages 13003–13051, Toronto,			
750	Canada. Association for Computational Linguistics.			
751	Gemini Team, Rohan Anil, Sebastian Borgeaud,			
752	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,			
753	Radu Soricut, Johan Schalkwyk, Andrew M Dai,			
754	Anja Hauth, et al. 2023. Gemini: a family of			
755	highly capable multimodal models. <i>arXiv preprint</i>			
756	<i>arXiv:2312.11805</i> .			
757	Xin Tian, Yingzhan Lin, Mengfei Song, Siqi Bao, Fan			
758	Wang, Huang He, Shuqi Sun, and Hua Wu. 2022.			
	Q-tod: A query-driven task-oriented dialogue system.			
	<i>arXiv preprint arXiv:2210.07564</i> .			
	Fanqi Wan, Weizhou Shen, Ke Yang, Xiaojun Quan, and			
	Wei Bi. 2023. Multi-grained knowledge retrieval for			
	end-to-end task-oriented dialog.			
	In <i>Annual Meeting of the Association for Computational Linguistics</i> .			
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten			
	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,			
	et al. 2022. Chain-of-thought prompting elicits rea-			
	soning in large language models. <i>Advances in neural</i>			
	<i>information processing systems</i> , 35:24824–24837.			
	Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas			
	Muennighoff. 2023. C-pack: Packaged resources			
	to advance general chinese embedding.			
	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,			
	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen			
	Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen			
	Yang, Yushuo Chen, Z. Chen, Jinhao Jiang, Ruiyang			
	Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu,			
	Jiayun Nie, and Ji rong Wen. 2023. A survey of			
	large language models.			
	<i>ArXiv</i> , abs/2303.18223.			
	A Experimental Details of Constructing			
	FakeRest Dataset			
	We utilize Prompt 1, 2, 3 to allow user LLMs and			
	system LLMs to have multiple rounds of dialogues			
	to build the dataset. Based on the user need path,			
	we pass a user need as input to the user LLM, such			
	as “find the moderate british restaurant in the center			
	area.”, and propose the need using Prompt 1. Based			
	on the need path we can also query the entities			
	from the knowledge base that match the need. If			
	the query finds an entity that matches the need, the			
	Prompt is used to reply to the user. If no entity is			
	found that matches the need, then Prompt is used			
	to apologize. We have the system LLM start the			
	dialogue with “What can I do for you?”, but we			
	don’t save this sentence in the dialogue dataset.			
	We used gemini-1.5-pro (Team et al., 2023) to			
	generate the dialog dataset. To generate a more			
	diverse set of responses, we use a temperature of			
	2.0. other hyperparameters use default values.			
	B Experimental Details of Fine-tuning			
	Embedding Models			
	In the retrieval process, we are using bge-large-			
	v1.5-en (Xiao et al., 2023) as the base model. We			
	utilize the AdamW optimizer (Loshchilov and Hut-			
	ter, 2017) and the linear learning rate scheduler			
	with 0.1 warmup steps. We set the epoch to 2 and			
	the learning rate to 2e-5. We utilize cosine similar-			
	ity to compute the most relevant set of entities in			
	the knowledge base.			

Prompt 1 Prompt for user LLM

You are a user. Please respond to assistant based on the need of this round.

Dialogue Context:

Assistant: What can I do for you?

User: Can you recommend a good British restaurant in the centre of town?

Assistant: There are a few British restaurants in the centre of town, including Bistro Delights, Epicurean Emporium, and Tastebud Tempstress.

Need: find the moderate british restaurant in the centre area.

Prompt 2 Prompt for system LLM

You are the assistant. Based on the 'Entities', please respond a sentence to the user.

Dialogue Context:

Assistant: What can I do for you?

User: Can you recommend a good British restaurant in the centre area?

Entities:

{ 'name': 'bistro delights', 'food': 'british', 'area': 'centre' },

{ 'name': 'epicurean emporium', 'food': 'british', 'area': 'centre' },

{ 'name': 'tastebud tempstress', 'food': 'british', 'area': 'centre' }

Prompt 3 Prompt for system LLM without entities

You are the assistant. Please express sorry for not finding the restaurant meets the needs.

Dialogue Context:

Assistant: What can I do for you?

User: I'm looking for a European restaurant in the north area. Could you recommend one?

Assistant: I recommend either Spice Haven or Tropical Treats, both of which serve European cuisine in the north area.

User: Actually, I'm looking for something a little more upscale. Do you have any other suggestions for expensive European restaurants in the north area?
