

# KETOD: Knowledge-Enriched Task-Oriented Dialogue

Anonymous ACL submission

## Abstract

Existing studies in dialogue system research mostly treat task-oriented dialogue and chit-chat as separate domains. Towards building a human-like assistant that can converse naturally and seamlessly with users, it is important to build a dialogue system that conducts both types of conversations effectively. In this work, we investigate how task-oriented dialogue and knowledge-grounded chit-chat can be effectively integrated into a single model. To this end, we create a new dataset, KETOD (Knowledge-Enriched Task-Oriented Dialogue), where we naturally enrich task-oriented dialogues with chit-chat based on relevant entity knowledge. We also propose two new models, SimpleToDPlus and Combiner, for the proposed task. Experimental results on both automatic and human evaluations show that the proposed methods can significantly improve the performance in knowledge-enriched response generation while maintaining a competitive task-oriented dialog performance. We believe our new dataset will be a valuable resource for future studies. The code and the dataset will be made publicly available.

## 1 Introduction

With the surge of deep neural networks and language model pre-training these years, there have been significant improvements in building conversational artificial intelligence. One major type of dialogue being studied is task-oriented dialogue (TOD) (Wen et al., 2017a; Budzianowski et al., 2018; Rastogi et al., 2020; Hosseini-Asl et al., 2020a), where the system aims to collect user intents/goals to complete certain tasks. In the current TOD systems, the system responses are mostly concise and templated, as we generally focus on the success of task completion but not providing a natural and engaging conversational experience. The latter is the target of another kind of popularly studied dialogue - knowledge-grounded chit-chat (Ghazvininejad et al., 2018; Zhang et al., 2018;

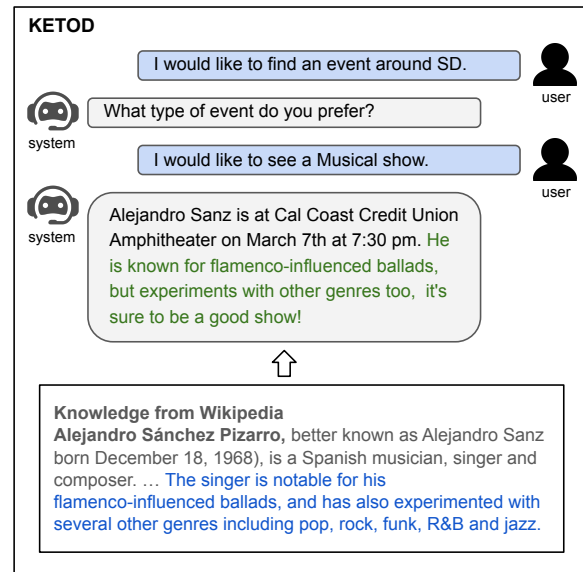


Figure 1: An example from the KETOD dataset: the green text is our enriched chit-chat based on the entity knowledge of *Alejandro Sanz* in the original TOD. Such knowledge-grounded chit-chat makes the dialogue more natural and engaging.

Tuan et al., 2019; Dinan et al., 2019). In addition to simple casual chit-chat, knowledge-grounded chit-chat further enables the conversation around certain background knowledge, relevant subjects or world knowledge.

Existing studies mostly focus on one specific type of dialogue and have gained remarkable progress over the last few years. However, the ultimate goal of Conversational AI is a human-like, unified system capable of conversing with the users naturally and seamlessly among all kinds of dialogues. Current TOD systems can hardly make interesting and engaging conversations only with templated functional responses. Few previous works like ACCENTOR (Sun et al., 2021) have studied the combination of TOD and chit-chat, but their chit-chat augmentation is largely limited to simple general responses like 'you're welcome', 'sounds good to me'. In this work, we propose to enrich TOD with knowledge-grounded chit-chat, as one

step further towards the ultimate goal of building a human-like, unified system (See Figure 1 for an example). We believe that the proposed knowledge-enriched TOD system can conduct more social, natural, and engaging conversations.

To this end, we propose a new dataset, KETOD (Knowledge-Enriched Task-Oriented Dialogue). In order to obtain natural and high-quality knowledge-grounded chit-chat, we design the dataset construction framework by augmenting existing TODs and using the relevant entity knowledge to make the chit-chat enrichment. Specifically, for a given task-oriented dialogue, 1) extracting the entities from the dialogue states and actions; 2) retrieving the knowledge associated with the entities from external knowledge sources; 3) asking the human annotators to enrich the system responses with chit-chat using the retrieved knowledge. We demonstrate that the knowledge-enriched dialogues constructed with the proposed framework are consistently preferred by human judges across all axes of engagingness, interestingness, knowledge, and humanness.

We propose two models for our task, and study the challenges and insights of our new dataset. The first model is an end-to-end language model that jointly learns and generates both the TOD results (dialogue states and actions) and the knowledge-enriched responses. The second model is a pipeline model that first generates the TOD results, then uses another response generation model to generate the knowledge-enriched responses. We run comprehensive experiments to demonstrate the improvement over to the baselines, and show that our models can generate better knowledge-enriched responses while maintaining competitive performance on the TOD task.

To summarize, we make the following major contributions:

- We propose the task of combining TOD and knowledge-grounded chit-chat.
- We construct a new large-scale dataset, KETOD, with high-quality, manually annotated dialogue responses enriched with knowledge-grounded chit-chat. We will release the dataset upon acceptance of the paper.
- We propose two models for our dataset, and carry comprehensive experiments to study the challenges and insights. We believe our

dataset should be a valuable resource for building a human-like conversational assistant.

## 2 Related Work

**Task-oriented dialogue.** As the major application in current industry dialogue systems, task-oriented dialogue (TOD) has been one of the most popular types of dialogue in the research community. There have been many works on building each component of the TOD system, such as dialogue state tracking, action prediction, and response generation (Wen et al., 2015, 2017b; Mrksic et al., 2017; Zhong et al., 2018; Eric et al., 2020; Liu et al., 2018; Peng et al., 2017; Serban et al., 2017; Zhou et al., 2017). Later works begin to investigate building end-to-end systems (Bordes et al., 2017; Liu et al., 2018, 2017; Xu et al., 2020). Following the trend of building unified pre-trained models in NLP research, most recent works on TOD also apply such language model pre-training style methods on building end-to-end systems (Hosseini-Asl et al., 2020a; Peng et al., 2020; Su et al., 2021). Such methods have achieved top performances on various datasets.

Popular datasets in TOD include the DSTC challenge series (Williams et al., 2016), MultiWOZ (Budzianowski et al., 2018), SGD (Rastogi et al., 2020), etc. As the primary goal of TOD is the successful completion of the functional tasks, the system responses mostly tend to be concise and templated. In this work, we believe it’s the ultimate goal of dialogue research to build a unified dialogue assistant to naturally and seamlessly converse with all kinds of dialogues. And we take the next step of combining TOD with knowledge-grounded chit-chat by enriching the functional responses with external knowledge.

**Chit-chat dialogue.** Another type of popular studied dialogue is chit-chat, with the goal of making a natural and engaging conversation. Apart from the ‘pure’ simple chit-chat that mostly covers plain and general responses, more works focus on knowledge groundings to achieve better specificity and engagingness, such as using user profiles (Zhang et al., 2018), social media contexts (Sordoni et al., 2015), or knowledge graphs (Tuan et al., 2019; Moon et al., 2019), etc. In this work, our enriched chit-chat is grounded on open-domain knowledge, similar as the WOW dataset (Dinan et al., 2019), where the system converses with the users about certain topics involving entity knowledge in an open-ended setting. In contrast, WOW specifically focuses on

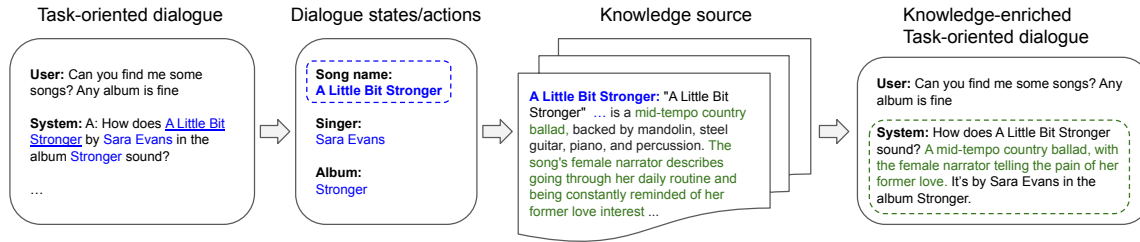


Figure 2: The pipeline of dataset construction: for each task-oriented dialogue, we first extract all the entities from the dialogue states and actions. Then we retrieve the knowledge associated with each entity from external knowledge sources (Wikipedia). At last, we ask human annotators to enrich the TOD system responses with chit-chat grounded on the retrieved knowledge.

knowledge-grounded chit-chat, while our dataset combines TOD and such chit-chat.

**Combination of task-oriented dialogue and chit-chat.** There are few works on combining TOD with chit-chat. ACCENTOR (Sun et al., 2021) proposes to combine TOD with chit-chat by prepending or appending chit-chat to the TOD system responses. But their chit-chat is mostly general responses like 'sounds good!', 'you're welcome'. In contrast, our work proposes to combine TOD with knowledge-grounded chit-chat. FusedChat (Young et al., 2021) proposes to insert chit-chat turns into TOD as well as re-writing TOD turns, but their chit-chat is still mostly general responses or based on commonsense knowledge, without much knowledge groundings.

### 3 The KETOD Dataset

#### 3.1 Dataset Construction

In this section, we describe our framework to construct the KETOD dataset. We start from existing TOD datasets and employ human annotators to augment the functional system responses with knowledge-grounded chit-chat. The proposed approach is demonstrated to give natural, contextual-relevant knowledge enrichment, and meanwhile easy to scale to different datasets. Figure 2 gives an overview of the dataset construction pipeline.

**Data preparation.** We build upon the SGD dataset (Rastogi et al., 2020), with task-oriented dialogues spanning 16 domains, such as *Restaurant*, *Wheather*, etc. Given each TOD, to obtain the knowledge relevant to the dialogue context, we first extract all the entities from the dialogue states and actions. We exclude the domains *Alarm*, *Banks*, and *Payment* as there are mostly no entities involved in these domains; Also, to simplify the human annotation process in the next step, we remove the dialogues with over 10 entities involved.

**Knowledge retrieval.** For each entity, we use the concatenation of the domain name and entity name as the query to retrieve Wikipedia articles. We use the DrQA retriever (Chen et al., 2017) to retrieve the top 2 Wikipedia articles and take the first 2 paragraphs of each article as the knowledge candidates associated with each entity. Then we break the retrieved articles into sentences, with each sentence as one knowledge snippet.

**Response enrichment.** In this step, we employ human annotators<sup>1</sup> to enrich the system responses in the original TOD based on the dialogue context and the retrieved knowledge. For each TOD, we present to the annotators the full dialogue, as well as all the knowledge snippets associated with the entities in the dialogue. The annotators can click on each entity name to see the associated knowledge snippets in an expanded textbox. See Appendix A for our annotation interface.

The annotation process is as follows: 1) Read the full dialog first to have an overall story in mind, as well as the relevant knowledge snippets, then to decide how many turns to enrich with chit-chat and which turn(s) to enrich; If there is no way to make a natural chit-chat enrichment, skip the example. 2) After deciding the turn(s) to enrich with the chit-chat, select the knowledge snippets used to make the enrichment (at most 3 snippets for each turn); 3) Rewrite the system response to enrich with chit-chat grounded on the selected knowledge snippets; The functional information in the original response should be maintained, while may be rephrased to make the enriched response more natural.

To ensure the dataset quality, we first interview the annotators to select the appropriate hires through a few test examples. Then we launch a training session for all the annotators to learn the

<sup>1</sup>Our annotators were hired as full-time employees through a leading annotation services vendor, and were paid in accordance with a fair wage rate.

|  |        |
|--|--------|
| Dialogues                              | 5,324  |
| Vocabulary                             | 27k    |
| All turns                              | 52,063 |
| Turns enriched with chit-chat          | 6,302  |
| All entities                           | 4,639  |
| All knowledge snippets                 | 33,761 |
| Avg. # turns per dialogue              | 9.78   |
| Avg. # tokens in enriched responses    | 28.07  |
| Avg. # entities per dialogue           | 4.98   |
| Avg. # knowledge snippets per dialogue | 70.50  |

Table 1: General statistics of KETOD.

task and the annotation interface. We launch the official batches after the annotators can well-master the task. During annotation, we specifically emphasize the contextualization of the knowledge-grounded chit-chat - the enrichment should be contextualized closely on the dialogue context, but not a plain restatement of the knowledge snippets.

### 3.2 Dataset Statistics and Analysis

We end up with 5,324 dialogues with enriched system responses. We make the split of 4,247/545/532 as the train/dev/test set. Table 1 shows the statistics of the KETOD dataset. Around 12.1% of the turns (which indicates mostly 1 or 2 turns in one dialogue) are enriched with knowledge-grounded chit-chat. This intuitively complies with our goal of making the whole dialogue natural and engaging, since too frequent chit-chat may result in redundancy and unnaturalness.

**Quality assessment of the annotation.** During the annotation process, around 12% of the dialogues cannot be enriched with any turns and thus discarded. To assess the quality of the annotation, we sample 5% of the annotated dialogues and distribute them to linguistics to check: 1) If the chit-chat enrichment is relevant and natural; 2) If the knowledge snippets are accurately selected corresponding to the enrichment. We end up with a correct rate of 87.0%.

**Justification of the chit-chat enrichment.** To demonstrate that our proposed knowledge-enriched TOD can be more natural and engaging, we conduct human evaluations to compare KETOD dialogues and their corresponding original TOD dialogues without chit-chat enrichment (SGD). We follow (Li et al., 2019) to make pairwise comparisons of the full dialogues over the following four axes: engagingness, interestingness, knowledge, and humanness. The results in Figure 3 show the superiority of KETOD over all axes.

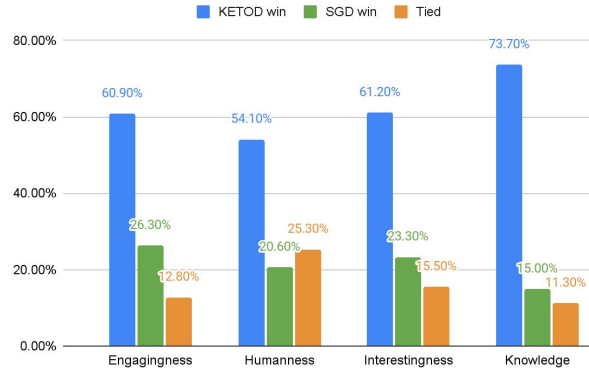


Figure 3: Results of pairwise comparison of KETOD vs SGD.

## 4 Approaches

In this section, we will describe the proposed two models for the KETOD dataset.

### 4.1 Overview and Formulations

For each dialogue turn, denote the dialogue context (history) as  $C$ , belief states as  $B$ , database search results as  $D$ , actions as  $A$ , the knowledge snippets used for chit-chat enrichment as  $K$ , the response as  $T$ . Then we formulate the problem as: given the dialogue context  $C$  and a knowledge source (Wikipedia in this dataset), the target is to generate the belief states  $B$ , actions  $A$ , and the response  $T$ , which may be enriched with chit-chat grounded on the knowledge based on the context. The goal of the optimization on KETOD is two-folded: 1) Optimizing the generation of knowledge-enriched responses; 2) Maintaining the task performances;

In this work, we propose the following modeling framework on KETOD: 1) given the dialogue context, generate the belief states and actions; 2) extract the entities in the belief states and actions, then use these entities to retrieve knowledge candidates (similar as in the dataset construction process); 3) conditioned on the dialogue context, use a knowledge selection model to select knowledge snippets from the knowledge candidates retrieved; 4) generate the knowledge-enriched response conditioned on both the dialogue context and the selected knowledge snippets.

Based on the above general framework, we propose two architectural approaches, **SimpleToD-Plus** and **Combiner**, respectively in §4.3 and §4.4.

### 4.2 Knowledge Selection

After the generation of belief states and actions, we retrieve the knowledge snippet candidates from Wikipedia using the entities in the belief states

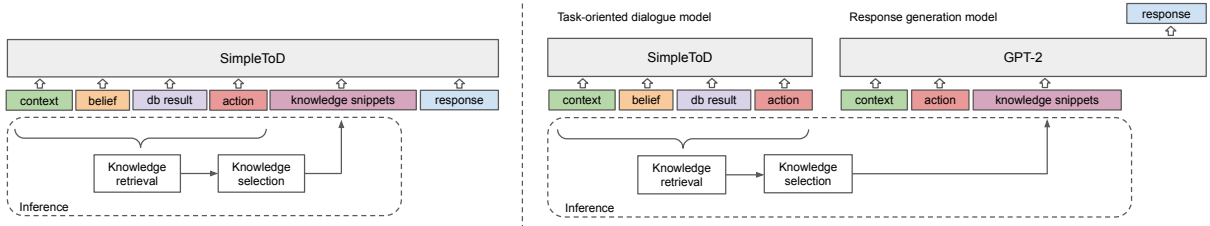


Figure 4: Illustration of the models. *Left*: the SimpleToDPlus model; *Right*: the Combiner model;

and actions. The average number of knowledge snippets candidates retrieved for each dialogue is around 70. It is impractical to input all of them into the models due to the large amount. As we have the annotation for the ground truth knowledge snippets used for each chit-chat enrichment, we train a knowledge selection model to select the top knowledge snippets most appropriate for chit-chat enrichment. Specifically, we concatenate the dialogue context with each knowledge snippet as the input. Then we use BERT (Devlin et al., 2019) to train a simple classifier to rank all the knowledge snippets candidates. We take the top 3 ones as the knowledge selection results. We use the same knowledge selection model for both architectures.

### 4.3 SimpleToDPlus

SimpleToD (Hosseini-Asl et al., 2020b) is a recent popular approach on TOD, which uses one single language model to sequentially generate the belief states, actions, and responses. It has achieved strong performances in all the above functional tasks. In this work, we propose its extension, **SimpleToDPlus**, to generate knowledge-enriched responses for TOD. The left part of Figure 4 shows the overview of SimpleToDPlus.

We formulate the training sequence as:

$$[C, B, D, A, K, \langle \text{chitchat} \rangle, T] \quad (1)$$

Where  $\langle \text{chitchat} \rangle$  is a tag to indicate the decision of whether to enrich the response with knowledge grounded chit-chat or not. If the response is not enriched, we insert the tag  $\langle \text{nochitchat} \rangle$ . Since the number of the gold knowledge snippets varies from 1 to 3 (as in the dataset construction), to be compatible with inference time, here we first run the knowledge selection model on all training instances. Then we construct the knowledge snippets  $K$  as the merge of the gold knowledge snippets and the knowledge selection model results, truncated to 3 ones. If the response is not enriched with chit-chat, i.e., no gold knowledge snippets, we still

put 3 snippets from the knowledge selection model ranking results here during training.

In the inference time, we first sequentially generate the belief states and actions. Then we extract the entities from the generated belief states and actions, and apply the same process of knowledge retrieval as in dataset construction. Next, we run the knowledge selection model on the retrieved knowledge candidates and take the top 3 knowledge snippets as the model input followed by the generated actions. At last, the model generates the decision to make chit-chat enrichment or not, followed by the final response.

Since the knowledge-enriched response is conditioned on the entity knowledge from the belief states and actions, we need to directly include the entities in the actions and responses during generation, instead of generating a delexicalized result first and then doing lexicalizing in the post-process as in the original SimpleToD. To simplify, we use the oracle database search results for all the experiments in this work.

### 4.4 Combiner

SimpleToDPlus models all the generations in an end-to-end manner. In **Combiner**, we use a pipeline of a TOD model followed by a response generation model to separate the TOD part (belief states, actions) with the generation of knowledge-enriched responses. The goal is to study whether an independent model can better learn each task with less interference from the other. The overview of the architecture is shown on the right of Figure 4.

For the TOD model, we use SimpleToD to generate the belief states and actions, with the training sequence as:

$$[C, B, D, A] \quad (2)$$

In the experiments, we find that including the knowledge-enriched responses here during training degrades the task performance, indicating the disturbance from the ungrounded knowledge in the responses.

| Models        | Joint GA    | Avg GA      | Act-Slot F1 | BLEU-4 <sub>aug</sub> | BLEU-4 <sub>orig</sub> | BLEU-4 <sub>all</sub> |
|---------------|-------------|-------------|-------------|-----------------------|------------------------|-----------------------|
| SimpleToD-ref | 27.6        | 54.2        | 67.6        | -                     | -                      | -                     |
| SimpleToD     | 23.7        | 50.1        | 62.7        | 4.8                   | 10.7                   | 10.0                  |
| SimpleToDPlus | <b>28.6</b> | <b>52.2</b> | <b>66.9</b> | 6.3                   | <b>11.7</b>            | <b>11.0</b>           |
| Combiner      | 24.5        | 51.5        | 64.5        | <b>6.5</b>            | 9.9                    | 9.5                   |

Table 2: Main experiment results: Both SimpleToDPlus and Combiner outperform the baseline. Overall SimpleToDPlus obtains better response generation performance while maintaining competitive TOD performance.

For the response generation model, we use GPT-2 (Radford et al., 2019) with the concatenation of the dialogue context, actions, and the knowledge snippets as the prompt:

$$T = \text{GPT-2}(C, A, K) \quad (3)$$

We use the same way of constructing the merged knowledge snippets during training, and the same process of knowledge retrieval and selection during inference as in SimpleToDPlus.

## 5 Experimental Results

**Baseline model.** We use the original SimpleToD (Hosseini-Asl et al., 2020b) as our baseline model, i.e., with the training sequence as  $[C, B, D, A, T]$ , without the injection of knowledge snippets. Therefore here, the knowledge-grounded chit-chat in the responses  $T$  do not have any knowledge groundings - we aim to show the necessity of knowledge grounding for our task, as well as the effectiveness of our proposed models to incorporate knowledge.

**Experimental setups and evaluations.** For all models, we use the Adam optimizer (Kingma and Ba, 2015). Check Appendix B for more details of training and parameter settings. For the TOD performances, we evaluate the belief states with joint goal accuracy (Joint GA) and average goal accuracy (Avg GA), and the actions with act-slot F1, same as (Sun et al., 2021). For the automatic evaluations of response generation, we use three BLEU-4 scores: BLEU-4<sub>aug</sub> for evaluating the responses enriched with knowledge; BLEU-4<sub>orig</sub> for evaluating the responses not enriched with knowledge; BLEU-4<sub>all</sub> for evaluating all responses;

### 5.1 Main Results

**Performance on response generation.** Table 2 shows our main experiment results. For the performances on response generation, we can see that both of our proposed models, SimpleToDPlus and Combiner, improve on the knowledge-enriched response generation (BLEU-4<sub>aug</sub>) over

the SimpleToD baseline. Since in the baseline, we do not include the knowledge snippets in the input, the generated responses are mostly enriched with random knowledge or frequent knowledge in the training data. The improvements demonstrate the necessity of knowledge grounding and the effectiveness of the proposed knowledge enrichment methods. Combiner performs slightly better on knowledge-enriched responses than SimpleToDPlus but falls short on the responses without knowledge-enrichment (i.e., original TOD responses). This is partially because of its pipeline nature - a separated response generation module can better learn the knowledge enrichment without the disturbance of other tasks, but the error cascading from the generated actions degrades the performance of the TOD responses part.

**Performances on belief states and actions.** To better study how the knowledge enrichment affects the TOD performances, we first train SimpleToD on our dataset without the knowledge enrichment, i.e., replace all the knowledge-enriched responses with the original responses in SGD. We name it as SimpleToD-ref in Table 2, serving as a reference of the original TOD performances. The SimpleToD baseline gives largely degraded performances due to the disturbance from the ungrounded knowledge in the responses during training. Therefore in Combiner, we do not include the responses in the training sequences of the TOD model (specified in section 4.4), and obtain better scores. SimpleToDPlus achieves the best TOD performances, which are nearly competitive with SimpleToD-ref.

**Human evaluations.** In order to get the more comprehensive measure of the response generation performances, we conduct human evaluations for both dialogue-level pairwise comparison and turn-level factualness evaluation. For dialogue-level pairwise comparison, we randomly sample 200 dialogues from the test set and apply the same process as in dataset evaluation (3.2). For each model, we construct the full dialogue results by concatenat-

| Metrics         | SimpleToDPlus win (%) | Combiner win (%) | Tied (%) |
|-----------------|-----------------------|------------------|----------|
| Engagingness    | 47.8                  | 24.5             | 27.8     |
| Interestingness | 34.5                  | 19.0             | 46.5     |
| Knowledge       | 29.5                  | 26.3             | 44.3     |
| Humanness       | 43.3                  | 23.8             | 33.0     |

Table 3: Human evaluation of SimpleToDPlus vs. Combiner.

| Metrics         | SimpleToDPlus win (%) | Gold win (%) | Tied (%) |
|-----------------|-----------------------|--------------|----------|
| Engagingness    | 16.8                  | 60.5         | 22.8     |
| Interestingness | 12.0                  | 51.0         | 37.0     |
| Knowledge       | 14.5                  | 44.8         | 40.8     |
| Humanness       | 17.3                  | 58.0         | 24.8     |

Table 4: Human evaluation of SimpleToDPlus vs. Gold.

476 ing the generated response for each turn given the  
477 gold dialogue context. Table 3 shows the results of  
478 pairwise comparison between the SimpleToDPlus  
479 model and the Combiner model, demonstrating  
480 SimpleToDPlus is more performant. Table 4 shows  
481 the results of pairwise comparison between Simple-  
482 ToDPlus and the gold reference, indicating there  
483 is still a large room for further improvements. For  
484 turn-level factualness evaluation, we randomly sam-  
485 ple one turn with chit-chat enrichment from each  
486 dialogue, and present both the generated response  
487 and the selected knowledge snippets to the anno-  
488 tators. The annotators are asked to check whether  
489 the chit-chat in the responses are factually correct  
490 based on the knowledge snippets. SimpleToDPlus  
491 and Combiner obtain the factualness correct rate of  
492 64.2% and 66.1%, respectively. In summary, Com-  
493 biner achieves better factualness of knowledge en-  
494 richment since its independent response generation  
495 model can better focus on the learning of knowl-  
496 edge groundings. But its error cascading due to the  
497 pipeline nature may degrade the overall consistency  
498 and human-likeness of the generated dialogue.

499 As we have two optimization goals on KE-  
500 TOD 1) Optimizing the generation of knowledge-  
501 enriched responses; 2) Maintaining the task perfor-  
502 mances, we consider SimpleToDPlus as a better  
503 model regarding the overall performances. We will  
504 use the results of SimpleToDPlus for the ablations  
505 and other analyses in the rest of the experiments.

## 5.2 Ablations and Analysis

507 **Analysis of different inference stages.** There are  
508 several inference stages for this task - the TOD  
509 results (belief states and actions), the selection of  
510 knowledge snippets, and the final response genera-

|  | BLEU-4 <sub>aug</sub> | BLEU-4 <sub>all</sub> |
|--|-----------------------|-----------------------|
| <b>Given gold TOD results, decision, and knowledge</b> |                       |                       |
| SimpleToD  | 6.5                   | 13.1                  |
| SimpleToDPlus  | 9.7                   | 14.6                  |
| Combiner   | 14.6                  | 15.1                  |
| <b>Given gold TOD results</b>                          |                       |                       |
| SimpleToD  | 6.3                   | 12.8                  |
| SimpleToDPlus  | 7.4                   | 14.0                  |
| Combiner   | 9.6                   | 13.9                  |

Table 5: Analysis of different inference stages: we provide the models with gold results up to certain stages, and investigate the performances for the inferences on following stages.

|                  | BLEU-4 <sub>aug</sub> | BLEU-4 <sub>all</sub> | Knowledge selection recall (%) |
|------------------|-----------------------|-----------------------|--------------------------------|
| Gold             | 9.7                   | 14.6                  | 100.0                          |
| BERT selection   | 7.8                   | 14.4                  | 52.7                           |
| TF-IDF selection | 6.6                   | 13.7                  | 14.1                           |

Table 6: SimpleToDPlus response generation performance with varying knowledge selection strategies.

511 tion, where each stage is conditioned on previous  
512 results. Therefore the errors accumulate through  
513 all the stages leading to the final performances.  
514 Here we run another two sets of experiments to  
515 study such error accumulations and compare the  
516 two models. Specifically, first, we feed the models  
517 with the gold TOD results, chit-chat decisions, and  
518 knowledge snippets, to solely test the abilities to  
519 generate the knowledge-enriched responses; Sec-  
520 ond, we feed the models with the gold TOD results  
521 to test the following stages of knowledge selec-  
522 tion and the response generation. The results are  
523 shown in Table 5. Compared with the full inference  
524 results in Table 2, we can see that the Combiner  
525 model largely outperforms SimpleToDPlus if pro-  
526 vided with more gold results for previous stages.  
527 However, it gradually falls behind SimpleToDPlus  
528 when moving to fully end-to-end inference due to  
529 the error cascading of its pipeline nature.

### Importance of knowledge selection strategies.

530 To demonstrate the importance of the knowledge  
531 selection strategies (and their subsequent recall per-  
532 formance), we run SimpleToDPlus with 1) gold  
533 knowledge snippets; 2) predicted knowledge snip-  
534 pets (with BERT); 3) knowledge snippets selected  
535 by heuristics (we use TF-IDF matching between  
536 the current dialogue turn and the knowledge snip-  
537 pets). To eliminate the influences brought by other  
538 inference stages, we feed the model with gold TOD  
539 results (dialogue states and actions). The results  
540 are shown in Table 6. There exists a certain level  
541 of variance for knowledge selection, e.g., when rec-  
542 ommending a song for the user, you may talk about  
543

| Gold   | Knowledge snippets   | Predicted   | Knowledge snippets   |
|--|--|---|--|
| <p><b>User:</b> Sure, that sounds right. Could you help me find some songs?<br/> <b>System:</b> I found 10 songs for you. One song is Supercut by Lorde from the album Melodrama.</p> <p><b>User:</b> What other songs did you have? I like Battle Beast.<br/> <b>System:</b> They're quite a well-known power metal band. I found 1 song for you. It is called Eden by Battle Beast from the album No More Hollywood Endings?</p> | <p><b>Battle Beast</b><br/> Battle Beast is a Finnish heavy/power metal band established in Helsinki in 2005. The initial line-up consisted of guitarists...</p>                                       | <p><b>User:</b> Sure, that sounds right. Could you help me find some songs?<br/> <b>System:</b> I have 10 songs for you. Would you like to listen to Supercut by Lorde from the album Melodrama?</p> <p><b>User:</b> What other songs did you have? I like Battle Beast.<br/> <b>System:</b> I found 1 song by Battle Beast from the album No More Hollywood Endings. It is a metal song by Battle Beast.</p> | <p><b>Battle Beast</b><br/> Battle Beast is a Finnish heavy/power metal band established in Helsinki in 2005. The initial line-up consisted of guitarists...</p> |
| <p><b>User:</b> I would like to find an event around SD.<br/> <b>System:</b> What type of event do you prefer?</p> <p><b>User:</b> I would like to see a Musical shows.<br/> <b>System:</b> Alejandro Sanz is at Cal Coast Credit Union Amphitheater on March 7th at 7:30 pm. He is known for flamenco-influenced ballads, but experiments with other genres too, so it's sure to be a good show!</p>                              | <p><b>Alejandro Sánchez Pizarro...</b> The singer is notable for his flamenco-influenced ballads, and has also experimented with several other genres including pop, rock, funk, R&amp;B and jazz.</p> | <p><b>User:</b> I would like to find an event around SD.<br/> <b>System:</b> What type of event would you like? Music, Sports, or something else?</p> <p><b>User:</b> I would like to see a Musical shows.<br/> <b>System:</b> I have 10 events that you might like. One is Alejandro Sanz at Cal Coast Credit Union Amphitheater. He is a Spanish singer-songwriter and musician.</p>                        | <p><b>Alejandro Sánchez Pizarro,</b> is a Spanish musician, singer and composer...</p>   |

Figure 5: Case studies: two examples of comparing the generation from SimpleToDPlus (right) with the gold reference (left), together with the knowledge snippets selected. Overall our model can mostly generate reasonable knowledge enrichment, but still falls short on engagingness and consistency compared to the gold references.

|                    | BLEU-4 <sub>aug</sub> | BLEU-4 <sub>all</sub> | Enrichment decision<br>F1 (%) |
|--------------------|-----------------------|-----------------------|-------------------------------|
| Gold decision      | 9.7                   | 14.6                  | 100.0                         |
| Predicted decision | 8.0                   | 14.1                  | 58.7                          |

Table 7: SimpleToDPlus response generation performance using (1) the gold set of turns to enrich with chit-chat, and (2) the predicted set of turns.

|                       | All  | Hotels | Movies | Restaurant | Music |
|-----------------------|------|--------|--------|------------|-------|
| BLEU-4 <sub>aug</sub> | 6.3  | 7.1    | 5.2    | 5.1        | 7.7   |
| BLEU-4 <sub>all</sub> | 11.0 | 10.3   | 12.2   | 14.0       | 12.3  |

Table 8: Domain breakdown of SimpleToDPlus response generation performances.

its genre, its singer, or the album.

**Learning when to inject knowledge-enriched chit-chat.** In all models, we use the special token ‘<chitchat>’ and ‘<nochitchat>’ to indicate the decision to inject knowledge enrichment for the responses. To study the effect of the chit-chat injection decision-making accuracy on the overall dialogue tasks, we run SimpleToDPlus (1) with the ground-truth information of turns to enrich with chit-chat, and (2) with the predicted decisions, using the gold TOD results. Table 7 shows the performance gap, which highlights the importance of knowing when to inject knowledge-enriched chit-chat. While such decisions are conditioned on the dialogue history, e.g., we may tend to not enrich a turn if many of the previous turns are enriched to avoid redundancy, there also exists some variance. In a real system, we may consider specifying the turns to make the chit-chat enrichment instead of letting the model make the decision.

**Domain analysis.** We investigate the model performance for each domain in Table 8. We observe that the performance differences may depend on the variance of the enriched knowledge. Domains with larger variance on selected knowledge tend to have lower automatic scores. For example, in `Hotels` domain, mostly the chit-chat is about the locations since there are mostly location entities involved in this domain. But for the `restaurants` domain, the enriched knowledge can be about the food, the restaurant, as well as the location. The selected knowledge shows more diversity and variance.

We provide case studies in Figure 5 to compare the predicted results with the gold references.

## 6 Conclusion and Future Work

In this work, we propose to combine task-oriented dialogue with knowledge-grounded chit-chat as one step further towards building a unified, human-like conversational assistant. To this end, we construct a new dataset named KETOD, with manually composed knowledge-enriched system responses. We conduct comprehensive experiments on our new dataset to study the insights and challenges. In real-world conversations, there are many ways that different kinds of dialogues are fused together. In this work, we start from the more straightforward type of directly enriching the TOD responses with knowledge-grounded chit-chat. Potential future works could be on studying other fusion strategies, for example, interleaving TOD dialogue turns with chit-chat turns, chit-chat grounded on multiple types of knowledge sources, etc.



## 7 Ethical Considerations

**Data Access and Licensing.** We develop the KETOD dataset based on the publicly available SGD dataset<sup>2</sup> (Rastogi et al., 2020). The SGD dataset is publicly available under the CC-BY-SA-4.0 License.

**Dataset Collection Process and Conditions.** For the annotation of our KETOD dataset, linguistics for assessing data quality, and all the human evaluations, our annotators were hired as full-time employees through a leading annotation services vendor, and were paid in accordance with a fair wage rate.

## References

Antoine Bordes, Y-Lan Boureau, and Jason Weston. 2017. [Learning end-to-end goal-oriented dialog](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.

Pawel Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gasic. 2018. [Multiwoz - A large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 5016–5026. Association for Computational Linguistics.

Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1870–1879. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. [Wizard of wikipedia: Knowledge-powered conversational agents](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

<sup>2</sup><https://github.com/google-research-datasets/dstc8-schema-guided-dialogue>

Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Kumar Goyal, Peter Ku, and Dilek Hakkani-Tür. 2020. [Multiwoz 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking baselines](#). In *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, pages 422–428. European Language Resources Association.

Marjan Ghazvininejad, Chris Brockett, Ming-Wei Chang, Bill Dolan, Jianfeng Gao, Wen-tau Yih, and Michel Galley. 2018. [A knowledge-grounded neural conversation model](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18), the 30th Innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 5110–5117. AAAI Press.

Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020a. [A simple language model for task-oriented dialogue](#). *CoRR*, abs/2005.00796.

Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020b. [A simple language model for task-oriented dialogue](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

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*.

Margaret Li, Jason Weston, and Stephen Roller. 2019. [ACUTE-EVAL: improved dialogue evaluation with optimized questions and multi-turn comparisons](#). *CoRR*, abs/1909.03087.

Bing Liu, Gökhan Tür, Dilek Hakkani-Tür, Pararth Shah, and Larry P. Heck. 2017. [End-to-end optimization of task-oriented dialogue model with deep reinforcement learning](#). *CoRR*, abs/1711.10712.

Bing Liu, Gökhan Tür, Dilek Hakkani-Tür, Pararth Shah, and Larry P. Heck. 2018. [Dialogue learning with human teaching and feedback in end-to-end trainable task-oriented dialogue systems](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 2060–2069. Association for Computational Linguistics.

Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. [Opendialkg: Explainable conversational reasoning with attention-based walks over](#)



821 with exact log-likelihood optimization. *World Wide*  
822 *Web*, 23(3):1989–2002.

and batch size of 16. All the experiments are done  
using TESLA M40 GPU cards.

871  
872

823 Tom Young, Frank Z. Xing, Vlad Pandealea, Jinjie  
824 Ni, and Erik Cambria. 2021. [Fusing task-oriented](#)  
825 [and open-domain dialogues in conversational agents](#).  
826 *CoRR*, abs/2109.04137.

827 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur  
828 Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you](#)  
829 [have pets too?](#) In *Proceedings of the 56th Annual*  
830 *Meeting of the Association for Computational Lin-*  
831 *guistics, ACL 2018, Melbourne, Australia, July 15-*  
832 *20, 2018, Volume 1: Long Papers*, pages 2204–2213.  
833 Association for Computational Linguistics.  
834

835 Victor Zhong, Caiming Xiong, and Richard Socher.  
836 2018. [Global-locally self-attentive encoder for di-](#)  
837 [alogue state tracking](#). In *Proceedings of the 56th An-*  
838 *ual Meeting of the Association for Computational*  
839 *Linguistics, ACL 2018, Melbourne, Australia, July*  
840 *15-20, 2018, Volume 1: Long Papers*, pages 1458–  
841 1467. Association for Computational Linguistics.

842 Ganbin Zhou, Ping Luo, Rongyu Cao, Fen Lin,  
843 Bo Chen, and Qing He. 2017. [Mechanism-aware](#)  
844 [neural machine for dialogue response generation](#). In  
845 *Proceedings of the Thirty-First AAAI Conference on*  
846 *Artificial Intelligence, February 4-9, 2017, San Fran-*  
847 *cisco, California, USA*, pages 3400–3407. AAAI  
848 Press.

## 849 **Appendix A: Dataset Construction**

850 Figure 6 shows our annotation interface to add  
851 knowledge-grounded chit-chat to TOD. The left  
852 part shows the full dialogue, where the annotators  
853 can click and expand each turn to make the chit-  
854 chat enrichment. The right part shows all the enti-  
855 ties with the associated knowledge snippets. The  
856 annotators can click on each entity name to expand  
857 the textbox to see the knowledge snippets. We add  
858 index number to each knowledge snippet (shown  
859 in green brackets), and the annotators are asked to  
860 write down the indexes of the knowledge snippets  
861 they used for writing the knowledge grounded chit-  
862 chat. Figure 7 shows one example annotation turn  
863 using our interface.

## 864 **Appendix B: Model and Training Details**

865 All the implementations are based on the Hugging-  
866 face Transformers library<sup>3</sup>. For the knowledge se-  
867 lection model, we use BERT-base with learning  
868 rate of 3e-5 and batch size of 16. For the base-  
869 line SimpleToD model, SimpleToDPlus model, and  
870 Combiner model, we all use learning rate of 1e-4

<sup>3</sup><https://github.com/huggingface/transformers>

# Enrich Task-Oriented Dialog With Entity Knowledge

Instructions ▾

## Dialog 15\_00100

Pre-check: (required)  
Select one ▾

**Turn 1**  
User: Can you get me a hotel?  
System: Where are you going.  
[Rewrite knowledge-enriched response for turn 1](#)

**Turn 2**  
User: I need a room in Los Angeles.  
System: Try this 2 star one: AC Hotel by Marriott Beverly Hills.  
[Rewrite knowledge-enriched response for turn 2](#)

**Rewrite knowledge-enriched response**

**Knowledge snippets selected**

**Turn 3**  
User: Where exactly are they? How much will this set me back?  
System: \$120 a night. Go to 6399 Wilshire Boulevard.  
[Rewrite knowledge-enriched response for turn 3](#)

**Turn 4**  
User: Alright. Thanks.  
System: Do you want me to got ahead with the reservation?  
[Rewrite knowledge-enriched response for turn 4](#)

**Turn 5**  
User: No. Get me a full-size car there until March 12th.  
System: From March 6th? Time?  
[Rewrite knowledge-enriched response for turn 5](#)

## Knowledge Snippets

About: Los Angeles

### History of Los Angeles

[1] The written history of Los Angeles city and county began with a small Mexican town that changed little in the three decades after 1848, when California became part of the United States.  
[2] Much greater changes came from the completion of the Santa Fe railroad line from Chicago to Los Angeles in 1885.  
[3] Immigrants flooded in, especially white Protestants from the Midwest.  
[4] Los Angeles had a strong economic base in farming, oil, tourism, real estate and movies.  
[5] It grew rapidly with many suburban areas inside and outside the city limits.  
[6] Hollywood made the city world famous, and World War II brought new industry, especially high-tech aircraft construction.  
[7] Politically the city was moderately conservative, with a weak labor union sector.  
[8] Since the 1960s growth has slowed—and traffic delays have become famous.  
[9] LA was a pioneer in freeway development as the public transit system deteriorated.  
[10] New arrivals, especially from Mexico and Asia, have transformed the demographic base since the 1960s.  
[11] Old industries have declined, including farming, oil, military and aircraft, but tourism, entertainment and high tech remain strong.

### Los Angeles

[12] Los Angeles (, Spanish for "The Angels"), officially the City of Los Angeles and often known by its initials L.A., is the second-most populous city in the United States (after New York City), the most populous city in California and the county seat of Los Angeles County.  
[13] Situated in Southern California, Los Angeles is known for its mediterranean climate, ethnic diversity, sprawling metropolis, and as a major center of the American entertainment industry.  
[14] Los Angeles lies in a large coastal basin surrounded on three sides by mountains reaching up to and over .  
[15] Historically home to the Chumash and Tongva, Los Angeles was claimed by Juan Rodríguez Cabrillo for Spain in 1542 along with the rest of what would become Alta California.  
[16] The city was officially founded on September 4, 1781, by Spanish governor Felipe de Neve.

Figure 6: Our annotation interface example 1.

Dialog 15\_00100

Pre-check: (required)  
Valid example ▾

**Turn 1**  
User: Can you get me a hotel?  
System: Where are you going.  
[Rewrite knowledge-enriched response for turn 1](#)

**Turn 2**  
User: I need a room in Los Angeles.  
System: Try this 2 star one: AC Hotel by Marriott Beverly Hills.  
[Rewrite knowledge-enriched response for turn 2](#)

**Rewrite knowledge-enriched response**

**Knowledge snippets selected**

**Turn 3**  
User: Where exactly are they? How much will this set me back?  
System: \$120 a night. Go to 6399 Wilshire Boulevard.

## Knowledge Snippets

About: Los Angeles

About: Ac Hotel By Marriott Beverly Hills

### Beverly Hills, California

[28] Beverly Hills is a city in Los Angeles County, California, United States, surrounded by the cities of Los Angeles and West Hollywood.  
[29] Originally a Spanish ranch where lima beans were grown, Beverly Hills was incorporated in 1914 by a group of investors who had failed to find oil, but found water instead and eventually decided to develop it into a town.  
[30] By 2013, its population had grown to 34,658.  
[31] Sometimes referred to as "90210", one of its primary ZIP codes, it was home to many actors and celebrities throughout the 20th century.  
[32] The city includes the Rodeo Drive shopping district and the Beverly Hills Oil Field.  
[33] Gaspar de Portolá arrived in the area that would become Beverly Hills on August 3, 1769, travelling along native trails which followed the present-day route of Wilshire Boulevard.

### The Beverly Hills Hotel

[34] The Beverly Hills Hotel, also called "The Beverly Hills Hotel and Bungalows", is located on Sunset Boulevard in Beverly Hills, California.  
[35] One of the world's best-known hotels, it is closely associated with Hollywood film stars, rock stars and celebrities.  
[36] The hotel has 208 guest rooms and suites, and 23 bungalows, each designed in the peachy pink and green colors which are a trademark of the hotel.

Figure 7: Our annotation interface example 2.