
Supplementary Materials

1 A Appendix

2 A.1 Construction & Schema Details

3 A.1.1 Conversation Details

4 To make spoken conversations that close to the real scenarios, we change the following interaction
5 pattern in MultiWOZ. In SpokenWOZ, once the user’s booking is successful, the agent will provide
6 the entity booked and ask for the user’s profile information, rather than providing a reference code in
7 MultiWOZ. Profile information including name, ID, email, license plate number, and phone. We will
8 explain in detail when agents will actively collect profile information from users.

9 **Name.** When a user makes a successful
10 hotel and restaurant reservation, the
11 agent will request the user’s name as the
12 reserved information. The user’s name
13 is randomly generated by the script¹.

14 **ID number.** When a user books a train,
15 the agent will ask for the user’s ID number
16 as registration information, which is
17 a randomly generated 16-digit string.

18 **Email.** When the user completes the
19 hotel or restaurant reservation, the agent
20 will ask the user if she/he wants to receive
21 the order via email. If the user
22 agrees to receive the order, the agent will
23 request the user’s email. The mailbox
24 number consists of the first letter of the
25 user’s first name plus the user’s last name, plus four randomly generated characters, and randomly
26 choose one of “@gmail.com”, “@yahoo.com”, “@outlook.com”, “@hotmail.com” as the suffix.

27 **License plate number.** When a user reserves a parking space at a hotel, attraction, or restaurant, the
28 agent will request the user’s license plate number. The license plate number is a string of 7 random
29 characters, the first two are letters, the middle two are numbers, and the last three are letters.

30 **Phone number.** When a user successfully books a taxi, the agent will request the user’s phone
31 number to contact the taxi driver, which is a randomly generated 10-digit string. In another case,
32 when users inquire about police station information, the agent will also ask for the user’s phone
33 number as a contact number.

¹names: <https://github.com/treyhunner/names>

Table 1: The 36 slots are tracked in the dialogue state.

<i>attraction</i>	area / name / type
<i>hospital</i>	department
<i>hotel</i>	area / bookday / bookpeople / bookstay / internet / name / parking / pricerange / stars / type
<i>restaurant</i>	booktime / bookday / bookpeople / area / food / name / pricerange
<i>taxi</i>	arriveby / departure / leaveat / destination
<i>train</i>	arriveby / departure / destination / leaveat / bookpeople / day
<i>profile</i>	license plate number / name / ID / email / phone

34 A.1.2 Slot Details

35 The following 36 slots are tracked in the dia-
36 logue state shown in Table 1. We also list the
37 reasoning slot in Table 2. To control the number
38 of cases where the value needs to be reasoned
39 about in the reasoning slots, we require partic-
40 ipants to implicitly express the values specified
41 in the task goal. 20% of the reasoning slot val-
42 ues will be automatically marked as requiring
43 implicit expression in the conversation. Mean-
44 while, co-reference annotation is already present
45 in SpokenWOZ. Instead of annotating pronouns,
46 we directly annotate the appropriate value in the
47 corresponding slot.

48 A.1.3 Speaker Origins Details

49 Considering the different laws of data access of
50 the different countries, we chose Canada, Singa-
51 pore, China, South Africa to collect the audio data, which will enable us to open source the audio data.
52 We also found that the cost of collecting audio data from Canada and Singapore is about three times
53 that of collecting from South Africa. Therefore, within the same budget, we chose to collect more
54 audios from South Africa, and we believe that a larger data set would further prompt the research in
55 the community. The distribution of speaker origins are shown in Table 3.

56 A.1.4 Audio Details

57 Our audio files are two-track. One track rep-
58 represents the voice of the user and the other rep-
59 represents the voice of the agent. Meanwhile, the
60 sample rate of our audio files is 8000Hz. Each
61 dialogue corresponds to an audio file, and each
62 word is recorded in the text annotation corre-
63 sponding to the word context, start time and end
64 time. To avoid the problem of overlapping ut-
65 terances, we follow the rules below during the
66 collection: (i) prohibit the agent from using the backchannel to interrupt the user; (ii) when a user uses
67 a backchannel expression, the agent should respond to the backchannel correctly, rather than ignoring
68 it and continuing the previous utterance. Finally, the word error rate of ASR is 6.1%, calculated from
69 the manually modified agent utterances and the agent utterances recognized by ASR tool.

70 A.1.5 Data Splits

71 SpokenWOZ is split into 4200/500/1000 dia-
72 logues in order by train/dev/test. More details
73 can be found in Table 4. The results of experi-
74 ments are evaluated by the test set.

75 A.2 Experiment Details

76 A.2.1 DST Baselines

77 **BERT+TripPy** TripPy makes use of copy mechanisms to fill slots. A slot is filled by one of
78 three copy mechanisms, including: (1) span prediction: values are directly extracted from the user’s

Table 2: Reasoning slot in SpokenWOZ. The upper script indicates which domains it belongs to. *: universal, 1: restaurant, 2: hotel, 3: attraction, 4: taxi, 5: train, 6: hospital, 7: police, 8: profile.

<i>Temporal Reasoning</i>	leaveat ^{4,5} arriveby ^{4,5} booktime ¹ / day ⁵ bookday ^{1,2}
<i>Mathematical Reasoning</i>	bookpeople ^{1,2,5} bookstay ²
<i>Semantic Reasoning</i>	type ^{1,2,3} area ^{1,2,3} internet ² department ⁶ parking ²

Table 3: The origins diversity of SpokenWOZ. Participants come from four different countries to improve the diversity of spoken conversations.

Country	Dialogues	Percentage	People	Percentage
Canada	500	8.77%	60	24%
Singapore	500	8.77%	40	16%
China	2100	36.84%	30	12%
South Africa	2600	45.61%	120	48%

Table 4: Statistics of SpokenWOZ.

Dataset	Train	Dev	Test
Audio Hours	183	22	44
Dialogues	4,200	500	1000
Turns	149,126	18,384	35,564
Tokens	1,672,984	204,644	396,933
Avg. Turns	35.50	36.77	35.56
Avg. Tokens	11.21	11.13	11.16

79 utterances; (2) inform operations: a value may be copied from the system’s inform operations; (3)
80 slot copy: a value may be copied over from a different slot.

81 **SPACE+TripPy** We use SPACE to replace the original encoder BERT in TripPy. SPACE is a
82 semi-supervised pre-trained conversation model learning from large-scale dialogue corpora with
83 limited annotations, which can be effectively fine-tuned on different downstream dialogue tasks.

84 **SPACE+WavLM+TripPy** To use both speech and text data, we concatenate the embeddings from
85 SPACE and WavLM. Then we use a Transformer encoder as the fusion layer to allow the interaction
86 between the different modalities. Then we use the fused outputs as the representations in the TripPy.

87 **UBAR** UBAR is acquired by fine-tuning the GPT-2 on the sequence of the entire dialogue session
88 which is composed of user utterance, dialogue state, database result, system act, and system response
89 of every dialogue turn. During the inference time, it formulates DST as a sequence-to-sequence task.
90 It takes the current user utterance, dialogue history, and the previously predicted dialogue state as
91 input, and gets the dialogue state of the current user utterance.

92 **SPACE** SPACE is a semi-supervised pre-trained conversation model learning from large-scale
93 dialogue corpora, which is based on UniLM [2]. Such as UBAR, we use SPACE as pre-trained
94 language model to fine-tuning on the sequence of the entire dialogue session. During the inference,
95 give the history and user utterance, SPACE generates dialogue states by autoregression. SPACE uses
96 the same special token as UBAR to split user utterances, dialogue state, act and system response.

97 **SPACE+WavLM** To utilize the generation ability of the SPACE model and the speech-modal
98 information, we concatenate the user utterance embeddings from SPACE and user audio embeddings
99 from WavLM as new user-side inputs. During the inference, the model uses dual-modal inputs to
100 generate the state by autoregression.

101 **SPACE+WavLM_{aligned}** Using the annotations in SpokenWOZ, the text of a word can be aligned
102 with its audio segment. To further explore how to make full use of the speech information, we align
103 the token and audio segment of every word in user utterances. Then we add the text embeddings from
104 SPACE and the corresponding embeddings from WavLM as new user-side embeddings. During the
105 inference, the model uses dual-modal input to generate the dialogue state by autoregression.

106 **ChatGPT:** ChatGPT (gpt-3.5-turbo) is a conversational LLM that has been trained by reinforcement
107 learning and instruction tuning [5], demonstrating a surprising ability in completing conversations.

108 **InstructGPT₀₀₃:** InstructGPT₀₀₃ (text-davinci-003) [5] is a 175B LLM trained by reinforcement
109 learning with human feedback and instruction tuning.

110 A.2.2 Response Generation Baselines

111 **UBAR** UBAR fine-tunes the GPT-2 on the sequence of the entire dialogue, including user utterance,
112 dialogue state, database result, system act, and system response. During the inference, UBAR uses
113 the fine-tuned GPT-2 to generate responses given different inputs based on different task settings.

114 **GALAXY** GALAXY is a pre-trained dialogue model that explicitly learns dialogue policy from
115 limited labeled dialogues and large-scale unlabeled dialogue corpora via semi-supervised learning,
116 which is based on UniLM. GALAXY keeps the same input format as UBAR.

117 **SPACE** SPACE is a pre-trained dialogue model that benefiting from large-scale dialogue corpora
118 via multi-task learning, including dialog understanding module, dialog policy module and dialog
119 generation module. SPACE is based on UniLM. SPACE keeps the same input format as UBAR.

120 **SPACE+WavLM** We concatenate the user utterance embeddings from SPACE and user audio
121 embeddings from WavLM as new user-side inputs. During the inference, SPACE+WavLM uses
122 dual-modal input to generate dialogue state, act, and final system response by autoregression.

123 **SPACE+WavLM_{aligned}** SPACE+WavLM_{aligned} adds the text embeddings from SPACE and the
124 corresponding embeddings from WavLM as new user-side embeddings. During the inference, the
125 model uses dual-modal input to generate dialogue state, act, and final system response.

126 A.2.3 Hyperparameters





127 For text-modal methods, we use the code and hyperparameters provided by their respective papers.
128 For dual-modal methods, we use the same hyperparameters as text-modal methods. To the fair
129 comparison, we train all the baselines 10 epoch for DST, 25 epoch for response generation and use
130 the final epoch checkpoint to get the results on SpokenWOZ. The results we report are the average of
131 the results using five different seeds. We trained the baselines in NVIDIA A100 and V100.

132 A.3 Case Study

133 We will show the predicted cases to confirm our insights proposed in section Experiments.




134 **Supervised generative methods are helpful.** We give the comparison between the generative-
135 method UBAR and extractive-method BERT+TripPy. Extractive methods can not get the value if it
136 does not directly exist in utterance, such as the reasoning slot. Meanwhile, the generative-method
137 can be robust to ASR noise and modify the wrong word in the utterance to the right one.

Table 5: The Case shows that supervised generative methods are helpful.

: There is a train, the id is called tr8925, do you want to make a booking.
: Yes, please make a booking for me.
: Okay. How many people?
: Um let me think ah **it's me and I have six friends with me.**
(BERT+TripPy: Train-Bookpeople = none)
(UBAR: Train-Bookpeople = 7)


138 The value 7 of slot Bookpeople can be correctly predicted by UBAR. However, value 7 is not existing
139 which is predicted wrongly by BERT+TripPy.


Table 6: The Case shows that supervised generative methods are helpful.


: Uh, I'm looking for a particular hotel in cambridge.
: Okay. So may I know its name, please?
: Um let me check. I think the name is called **lavelle lodge.**
(BERT+TripPy: Hotel-Name = lavelle lodge)
(UBAR: Hotel-Name = lovell lodge)


140 In this case, the correct name of the hotel "lovel lodge" is not predicted correctly by BERT + TripPy,
141 even it extracts the correct span in utterance. We also find that UBAR can get the name correctly,
142 which shows that generative-method can learn the ability to correct errors from ASR.


Table 7: The Case shows that supervised generative methods are helpful.

: My id is **88716**.
 (BERT+TripPy: Profile-ID = None)
 (UBAR: Profile-ID = 88716)

: You can proceed.

: **46859**.
 (BERT+TripPy: Profile-ID = None)
 (UBAR: Profile-ID = 8871646859)






: I am still working with.

: **638141**.
 (BERT+TripPy: Profile-ID = None)
 (UBAR: Profile-ID = 8871646859638141)

143 The generative-method can learn the ability to concatenate the value segment from different turns,
 144 which can be hardly learned by the extractive-method BERT + TripPy.






145 **Dual-modal TOD models is what you need.** We give the case study and comparison between
 146 text-modal methods SPACE + TripPy, SPACE and dual-modal methods SPACE + TripPy + WavLM,
 147 SPACE + WavLM. Although the main experimental results reflect that speech information improves
 148 overall performance, we are more concerned with the performance of the ASR-sensitive slots.

Table 8: The Case shows that dual-modal TOD models is what you need.

: Good afternoon. Yes. Uh, could you please assist me looking for a particular restaurant, please.
: No problem. Do you have a name for a restaurant?
: Yes. um, it's called the **bangkok cutting**.
: **Bangkok city**, okay. Just give me a second. I'll look for it on the system.
: Okay, not a problem at all.
 (SPACE+TripPy: Restaurant-Name = bangkok cutting)
 (SPACE+WavLM+TripPy: Restaurant-Name = bangkok city)
 (SPACE: Restaurant-Name = bangkok city)
 (SPACE+WavLM: Restaurant-Name = bangkok city)

149 In this case, the value of the slot Name is not correctly predicted by SPACE+TripPy. We find that the
 150 span prediction copy mechanism is performed in SPACE+TripPy. However, SPACE+WavLM+TripPy
 151 performed the copy mechanism and copy the value from Inform Operations to get right value. This
 152 indicates that the speech information can be used to further improve performance.

Table 9: The Case shows that dual-modal TOD models is what you need.

: I want to find a particular hotel to rise. I remember his name, but I can't find the location.
: So may I have the name of the hotel, please.
: Uh. the name of the hotel is **work worth house**.
: Okay, so I found a hotel called **warkworth house** for you.
: Uh, yes, yes. That's the hotel. Thank you.
 (SPACE+TripPy: Hotel-Name = None)
 (SPACE+WavLM+TripPy: Hotel-Name = work worth house)
 (SPACE: Hotel-Name = worth house)
 (SPACE+WavLM: Hotel-Name = warkworth house)

153 In this case, SPACE can not predict the value in the right pattern, but dual-modal
 154 SPACE+WavLM+TripPy successfully predict it. This shows that speech modality can also help
 155 generative methods learn the correct pattern, even in the presence of ASR noise from user utterances.

156 **A.4 Statistics**

157 We give the distribution of domains in Table
 158 10. Meanwhile, the dataset distributions of di-
 159 alog length and turn length are shown in the
 160 following figures. We give the statistics of Spo-
 161 kenWOZ in Figure 1 and 2. Shown in Figure
 162 1, the length of dialogue history is concentrated
 163 above 30 turns. The excessive number of dia-
 164 logue turns also makes it difficult for the model
 165 to learn. We compared the multi-domain and
 166 single-domain dialogue in Figure 3, intuitively,
 167 the number of turns for multi-domain dialogue
 168 is larger than the number of turns for single-
 169 domain dialogue. In Figure 4, there is no signif-
 170 icant difference between the utterance lengths
 171 of the user and agent, because SpokenWOZ is
 172 constructed using the Human-to-Human schema.
 173 We also show the distribution of the dialogue
 174 acts and slots in Figure 5.

Table 10: The distribution of dataset domains.

Domains	Number
profile-restaurant-train	720
hotel-profile-train	702
attraction-hotel-profile-taxi	295
hotel-profile-restaurant-taxi	294
attraction-profile-train	291
profile-taxi	285
attraction-train	278
attraction-profile-restaurant-taxi	275
hotel-profile-restaurant	252
profile-restaurant	238
attraction	238
hotel-profile	237
attraction-hotel-profile	212
attraction-profile-restaurant	209
profile-train	193
train	149
hotel-train	148
restaurant-train	132
hotel	104
restaurant	102
attraction-restaurant	85
attraction-hotel	62
hospital	57
police-profile	56
attraction-profile	47
hotel-restaurant	39
Total	5700

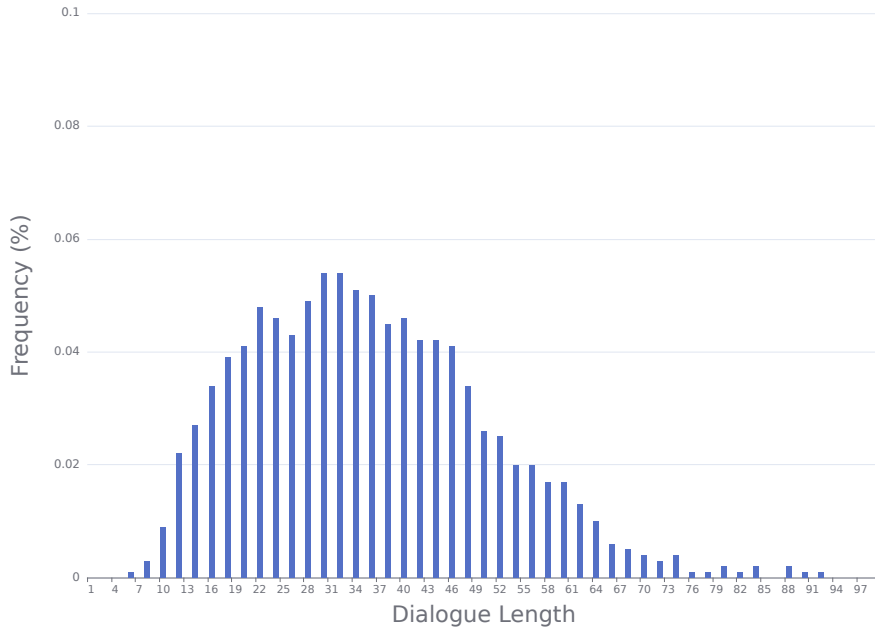


Figure 1: The distribution of the length of turn in SpokenWOZ.

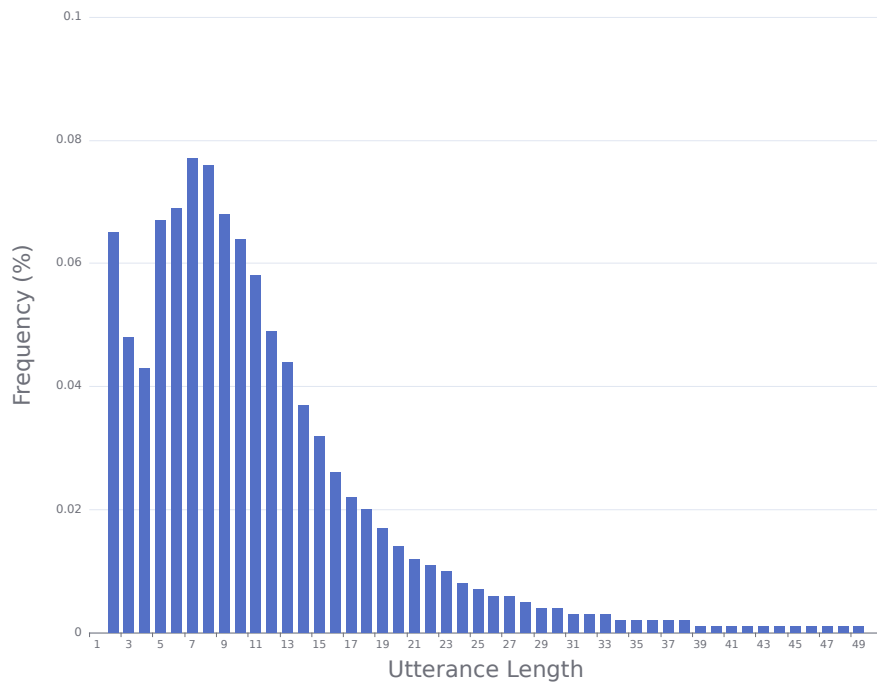


Figure 2: The distribution of the length of turn in SpokenWOZ.

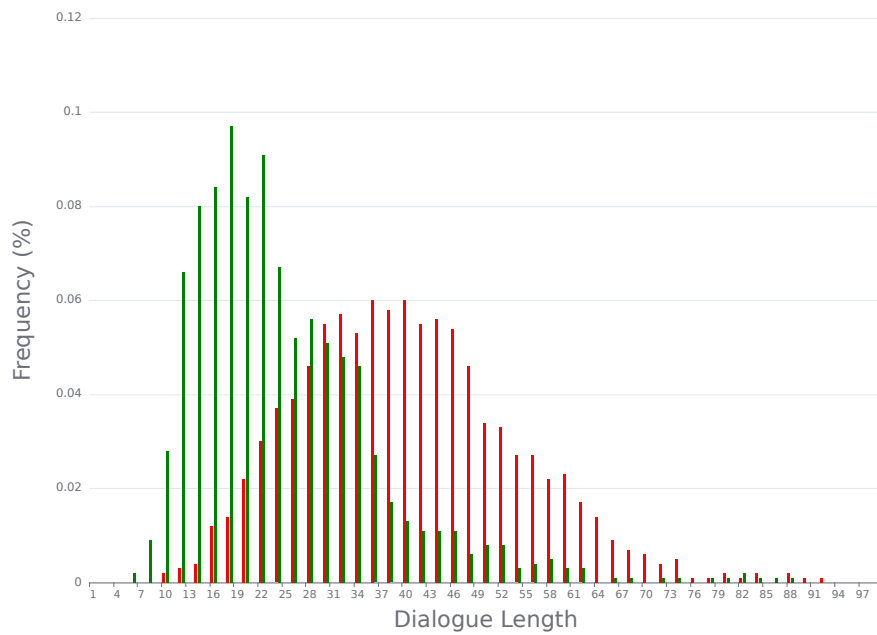


Figure 3: The distribution of the number of turns in two kinds of dialog in SpokenWOZ: **Multi-domain**, **Single-domain**.

180 **A.5 Heatmap of acts**

181 We show the act flow in SpokenWOZ in Figure 7. Given the user dialogue act, we present the
 182 frequency of agent act in heat map. As shown in Figure 6 and 7, SpokenWOZ not only contains more
 183 types of acts, but also contains more diverse act flow. It is more difficult for the model to predict the
 184 right action and give the right response.

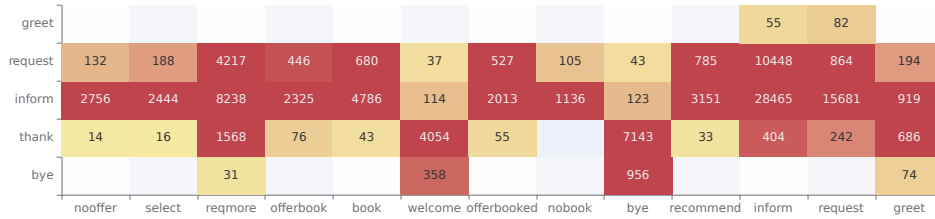


Figure 6: Heat map of agent acts in MultiWOZ. The heat map shows the frequency of the agent act (horizontal axis) after the given user act (vertical axis).

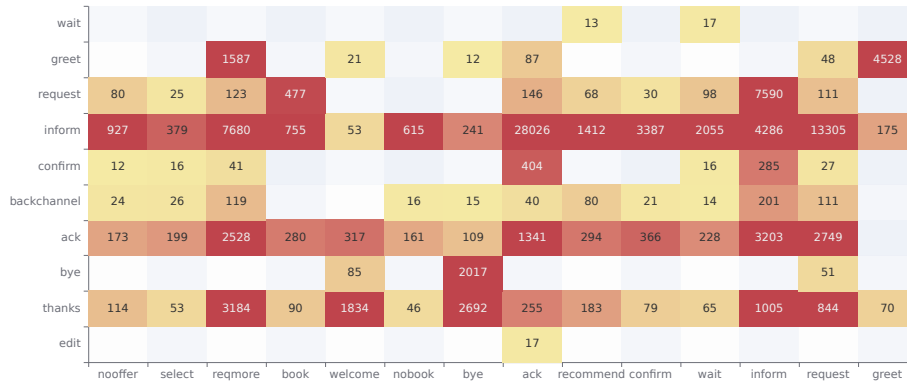


Figure 7: Heat map of agent acts in SpokenWOZ. The heat map shows the frequency of the agent act (horizontal axis) after the given user act (vertical axis).

185 **A.6 Dialogue Example**

186 A dialogue example can be found in Figure 8.



Figure 8: A dialogue example from SpokenWOZ.

187 **A.7 Online Database**

We give the interface of our built database for participants in Figure 9.

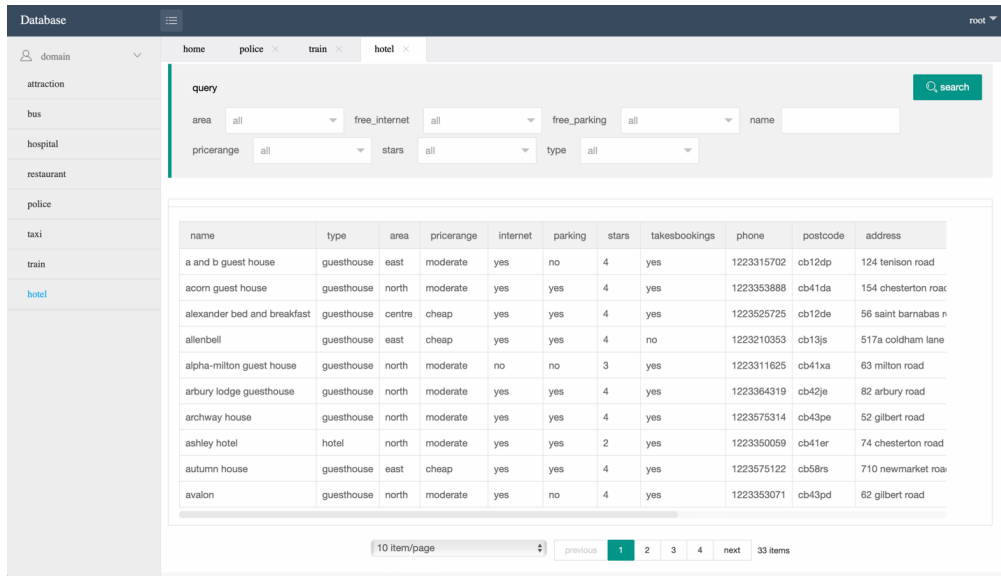


Figure 9: The built online database in SpokenWOZ.

188

189 **A.8 Task Goal Example**

190 We give an example of our task goal for participants in Table 11.

Table 11: An example of task goal.

(Tips!) The important information is marked with <>. Message with (use veiled expression) means that an veiled expression is needed here.

(Attention!) Please ask customer service for information and try to solve your problem in the Shortest number of turns

(Background) Your name is <Misty Imbert>, your telephone is <9310729130>, your ID number is <8485147021469113>, your email is <MImbert06jn@outlook.com>, your car_number is <OL09ODU>

(Background) You are planning your trip in Cambridge

You are looking for a <place to stay>. The hotel <doesn't need to have free parking> and should <include free wifi>

The hotel should be in the <west> and should have <a star of 4>

If there is no such hotel, how about one that has <free parking> Once you find the <hotel> you want to book it for <4 people> and <2 nights> starting from <friday>

You are also looking for a <restaurant>. The restaurant should be in the <moderate> price range and should serve <indian> food

The restaurant should be <in the same area as the hotel>

Once you find the <restaurant> you want to book a table for <the same group of people> at <12:15>(use veiled expression)> on <the same day>

(Background) Once you have made a booking, you do not want to give out your email address for receiving orders.

191 A.9 Analysis of LLMs

192 A.9.1 Prompts for LLMs

193 We imitate the format of prompt form Hudecek et al. [3] and Bang et al. [1]. We list a prompt
194 example for DST task in Table 12.

Table 12: An example of a zero-shot version of the prompt used for DST.

Definition: Give the dialogue state of the last utterance in the following dialogue in JSON (for example: STATE: "hotel-parking": "yes", "hotel-type": "guest house") by using the following pre-defined slots and possible values:

- Slot Name: hotel-pricerange; Slot Description: price budget of the hotel; Possible values: ['expensive', 'cheap', 'moderate']
- Slot Name: hotel-type; Slot Description: type of the hotel; Possible values: ['guest house', 'hotel']
- Slot Name: hotel-parking; Slot Description: whether the hotel has parking; Possible values: ['no', 'yes']
- Slot Name: hotel-day; Slot Description: day of the hotel booking; Possible values: ['monday', 'tuesday', 'wednesday', 'thursday', 'friday', 'saturday', 'sunday']
- Slot Name: hotel-people; Slot Description: number of people booking the hotel; Possible values: ['1', '2', '3', '4', '5', '6', '7', '8']
- Slot Name: hotel-stay; Slot Description: length of stay at the hotel; Possible values: ['1', '2', '3', '4', '5', '6', '7', '8']
- Slot Name: hotel-internet; Slot Description: whether the hotel has the free internet; Possible values: ['no', 'yes']
- Slot Name: hotel-name; Slot Description: name of the hotel; Possible values: []
- Slot Name: hotel-area; Slot Description: area of the hotel; Possible values: ['centre', 'east', 'north', 'south', 'west']
- Slot Name: hotel-star; Slot Description: star of the hotel; Possible values: ['0', '1', '2', '3', '4', '5']
- Slot Name: train-arriveby; Slot Description: the arrival time of the train, 24-hour standard time, e.g. 06:00, 18:30; Possible values: []
- Slot Name: train-day; Slot Description: day of the train departure; Possible values: ['monday', 'tuesday', 'wednesday', 'thursday', 'friday', 'saturday', 'sunday']
- Slot Name: train-people; Slot Description: number of people travelling by train; Possible values: ['1', '2', '3', '4', '5', '6', '7', '8']
- Slot Name: train-leaveat; Slot Description: leaving time of the train, 24-hour standard time, e.g. 06:00, 18:30; Possible values: []
- Slot Name: train-destination; Slot Description: destination of the train; Possible values: ['birmingham new street', 'bishops stortford', 'broxbourne', 'cambridge', 'ely', 'kings lynn', 'leicester', 'london kings cross', 'london liverpool street', 'norwich', 'peterborough', 'stansted airport', 'stevenage']
- Slot Name: train-departure; Slot Description: departure of the train; Possible values: ['birmingham new street', 'bishops stortford', 'broxbourne', 'cambridge', 'ely', 'kings lynn', 'leicester', 'london kings cross', 'london liverpool street', 'norwich', 'peterborough', 'stansted airport', 'stevenage']
- Slot Name: attraction-area; Slot Description: area of the attraction; Possible values: ['centre', 'east', 'north', 'south', 'west']
- Slot Name: attraction-name; Slot Description: name of the attraction; Possible values: []
- Slot Name: attraction-type; Slot Description: type of the attraction; Possible values: ['architecture', 'boat', 'cinema', 'college', 'concerthall', 'entertainment', 'museum', 'multiple sports', 'nightclub', 'park', 'swimmingpool', 'theatre']
- Slot Name: restaurant-pricerange; Slot Description: price budget for the restaurant; Possible values: ['expensive', 'cheap', 'moderate']
- Slot Name: restaurant-area; Slot Description: area of the restaurant; Possible values: ['centre', 'east', 'north', 'south', 'west']
- Slot Name: restaurant-food; Slot Description: the cuisine of the restaurant; Possible values: []
- Slot Name: restaurant-name; Slot Description: name of the restaurant; Possible values: []
- Slot Name: restaurant-day; Slot Description: day of the restaurant booking; Possible values: ['monday', 'tuesday', 'wednesday', 'thursday', 'friday', 'saturday', 'sunday']
- Slot Name: restaurant-people; Slot Description: number of people for the restaurant booking; Possible values: ['1', '2', '3', '4', '5', '6', '7', '8']
- Slot Name: restaurant-time; Slot Description: time of the restaurant booking, 24-hour standard time, e.g. 06:00, 18:30; Possible values: []
- Slot Name: hospital-department; Slot Description: department of the hospital; Possible values: []
- Slot Name: taxi-leaveat; Slot Description: leaving time of taxi, 24-hour standard time, e.g. 06:00, 18:30; Possible values: []
- Slot Name: taxi-destination; Slot Description: destination of taxi; Possible values: []
- Slot Name: taxi-departure; Slot Description: departure location of taxi; Possible values: []
- Slot Name: taxi-arriveby; Slot Description: arrival time of taxi, 24-hour standard time, e.g. 06:00, 18:30; Possible values: []
- Slot Name: profile-name; Slot Description: the name of the user; Possible values: []
- Slot Name: profile-email; Slot Description: the email of the user; Possible values: []
- Slot Name: profile-idnumber; Slot Description: the idnumber of the user; Possible values: []
- Slot Name: profile-phonenum; Slot Description: the phonenum of the user; Possible values: []
- Slot Name: profile-platenumber; Slot Description: the platenumber of the user; Possible values: []

USER: Hello, I need some help.
SYSTEM: Okay, how can i help you.
USER: Yes, I'm looking for a train.

STATE:

195 For response generation tasks, we follow the prompt as Hudecek et al. [3]. We use the dialogue state
196 to query the database provided to get the entity. For the Policy Optimization task, we use the ground
197 truth dialogue state to query the provided database. For End-to-end Modeling, we use the predicted

198 dialogue state generated from ChatGPT or InstructGPT₀₀₃ to query the provided database. We use
 199 the contents of the database and the generated system response to match the strings and complete
 200 the de-lexicalization process. The de-lexicalized response will be used to calculate INFORM and
 201 SUCCESS. We list a prompt example for a response generation task in Table 13.

Table 13: An example of a zero-shot version of the prompt used for Response Generation.

Definition: Please continue the dialogue as a task-oriented dialogue system called SYSTEM. The answer of SYSTEM should follow the DATABASE provided in JSON format and answer the USER’s last utterance. SYSTEM can recommend and inform the contents in the DATABASE according to the utterance of the USER and return the name of the entity when it comes to restaurants, hotels and attractions, and the trainid when it comes to trains. But only when the USER requests information about an entity in the DATABASE, such as a phone number, should SYSTEM inform the corresponding content.

DATABASE:
 { "restaurant": { "address": "Cambridge Retail Park Newmarket Road Fen Ditton", "area": "east", "food": "italian",
 "name": "pizza hut fen ditton", "phone": "12233237370", "postcode": "cb58wr", "pricerange": "moderate" } }




DIALOGUE CONTENT:
 USER: Hello.
 SYSTEM: Hello! How can I assist you today? Are you looking for any specific information or recommendations?
 USER: I'm in east part of the city and I want to have some food. Uh, hope for the. to serve Italian food and with moderate price range.

SYSTEM:

202 A.9.2 Analysis on DST




203 As shown in the section Experiments, the performances of LLMs do not surpass that of supervised
 204 small models and show a noticeable gap compared to the supervised generative dual-modal baselines
 205 in the DST task. However, it is worth noting that LLMs slightly outperform the BERT+TripPy when
 206 cross-turn slots are not taken into account, showing the potential capabilities of LLMs. Meanwhile,
 207 we find that the main reason for the poor results of LLMs is that the hallucination phenomenon [4] is
 208 very serious, e.g., LLMs often generate additional results that do not fit the dialogue utterances. When
 209 we measure only the slots where the ground truth value is not “none”, we find that ChatGPT achieves
 210 JGA of 30.81 and InstructGPT₀₀₃ achieves JGA of 34.42. Compared to the standard JGA, ChatGPT
 211 and InstructGPT₀₀₃ improve their scores by 17.06 and 20.27 respectively. This indicates that LLMs
 212 generate a large number of erroneous results at slots that are not involved in the conversation, i.e., the
 213 hallucination phenomenon is very serious. We show a case in Table 16.

Table 14: The Case shows the hallucination phenomenon.

: Hello.
: Hello, how can I help?
: Yes, I'm looking for restaurant.
 (ChatGPT: Restaurant-Food = international)
 (Ground Truth: Restaurant-Food = none)

214 Meanwhile, due to the inability of LLMs to perceive the information of speech, LLMs tend to generate
 215 the value directly from user utterance. We show a case in Table 15.


Table 15: The Case shows that LLMs are sensitive to the noisy utterance.

: Uh, I'm looking for a particular hotel in cambridge.
: Okay. So may I know its name, please?
: Um let me check. I think the name is called **lavelle lodge**.
 (ChatGPT: Hotel-Name = lavelle lodge)
 (InstructGPT₀₀₃: Hotel-Name = lavelle lodge)
 (Ground Truth: Hotel-Name = lovell lodge)

216 A.9.3 Analysis on Response Generation

217 As introduced in the section Experiments, LLMs achieve comparable performances in Policy Opti-
218 mization task but poor performances in the End-to-end Modeling task. We find that the main reason
219 for the poor results in End-to-end Modeling is that the entity returned by the database does not meet
220 the user’s needs. Meanwhile, the poor performances of BLEU show that there is a big difference
221 between LLM’s response style and human response style. Meanwhile, LLMs may feel confused
222 about noisy utterances and generate a statement requesting clarification as shown in Table 16.

Table 16: The Case shows that LLMs are sensitive to the noisy utterance.

: Mm. I’m looking for a place to die (*In the audio it is actually “dinner”*).
ChatGPT: I’m sorry, I’m not sure I understand. Could you please rephrase your question?

223 B Limitations

224 Even though we tried hard to build a realistic spoken TOD benchmark for further studies, we
225 could not use the audio data from real conversations due to privacy concerns. This brings the
226 following limitations: (1) for “profile” domain, we use a designed script to generate random personal
227 information, which may not be realistic, and the number of slots also limits the further in scenario that
228 agents need to collect personal information from users; (2) as the limited ontology, our benchmark
229 should mainly be used for research instead of deployed in realistic applications.

230 C Ethics Statement

231 We construct SpokenWOZ, a task-oriented dialogue benchmark containing both audio data and text
232 data under the CC BY-NC 4.0 License. We state that we bear all responsibility in case of violation
233 of rights. We will subsequently host and maintain the dataset in the corresponding website, and
234 welcome other researchers to improve the quality of SpokenWOZ together. Then, We will detail our
235 ethical considerations for each part of our collection process:

236 **Ontology Consideration.** We inherited and expanded MultiWOZ’s ontology, which is open-source
237 and under the MIT License. We have used it in compliance with its terms of use.

238 **Privacy Concern.** In our dataset, we have designed the scenarios where an agent needs to
239 proactively collect user information, however, our user information is all generated by scripts in a
240 random manner, so there will not be any privacy leakage issues.

241 **Audio Collection.** We informed each participant that the collected audio data will be used as
242 public dataset for research. Participants who agree to participate in data collection will sign a contract
243 with us, and the ownership and use rights of their data belong to us. We will not disclose which
244 specific participant the audio came from. At the same time, we have reviewed the legal regulations
245 of four countries and regions, and will release the data in a legal manner, so this will not cause any
246 legal problems. The distribution of our data sources has been discussed in the Appendix A.1.3, the
247 diversity of SpokenWOZ makes our data unbiased. We pay \$30k for 249 hours audio. The average
248 cost per hour of audio is \$120.

249 **Dialogue Annotation.** During the annotation process, the annotators’ personal information is not
250 collected, which will not cause privacy leakage. At the same time, during the annotation process, we
251 signed a non-disclosure agreement with each annotator, therefore, the audio data will not be leaked
252 during the annotation process. We pay \$20k for 5,700 dialogues. The average cost per dialogue
253 annotation is \$3.5. After our statistics, an average of one hour can annotate 5 dialogues.

254 D Data Format

255 D.1 Audio Format

256 As detailed in Appendix A.1.4, audio files are two-track with a sample rate of 8000Hz. One track
257 represents the voice of the user and the other represents the voice of the agent. Each dialogue
258 corresponds to an audio file, and each word is recorded in the text annotation corresponding to the
259 word context, start time, and end time. We use the wav format to save our audio files. The file name
260 of the audio is consistent with the id of the dialogue, for example, the corresponding audio file for
261 MUL0032 is MUL0032.wav.

262 D.2 Text Format

263 Our text data is given in json format, and we take the same fields as popular MultiWOZ 2.2 [6]
264 to store the corresponding information, so researchers can easily use our data. In addition, we
265 additionally provide the data ontology and the database json files. There are 5,700 dialogues ranging
266 from single-domain to multi-domain in SpokenWOZ. The test sets contain 1k examples. Dialogues
267 with MUL in the name refers to multi-domain dialogues. Dialogues with SNG refers to single-domain
268 dialogues. Each dialogue consists of a goal, multiple user and system utterances, dialogue state,
269 dialogue act, corresponding audio and ASR transcription.

270 The dialogue goal for each dialogue is recorded in the "goal" field. The dialogue goal holds the fields
271 involved in the dialogue as well as the slots involved and the corresponding values.

272 The dialogue state for each dialogue is recorded in the "metadata" field in every turn the same as
273 MultiWOZ 2.2. The state have two sections: semi, book. Semi refers to slots from a particular
274 domain. Book refers to booking slots for a particular domain. The joint accuracy metrics includes
275 ALL slots.

276 The dialogue act for each dialogue is recorded in the "dialogue_act" and "span_info" field in every
277 turn:

```
278 {  
279   "$dialogue_id": {  
280     "log": {  
281       "$turn_id": {  
282         "dialogue_act": {  
283           "$act_name": [  
284             [  
285               "$slot_name",  
286               "$action_value"  
287             ]  
288           ]  
289         },  
290         "span_info": [  
291           [  
292             "$act_name",  
293             "$slot_name",  
294             "$action_value",  
295             "$start_character_index",  
296             "$exclusive_end_character_index"  
297           ]  
298         ]  
299       }  
300     }  
301   }
```

302
303 The ASR transcription for each dialogue is recorded in the "words" field in every turn.

```
304 {  
305   "$dialogue_id": {  
306     "log": {  
307       "$turn_id": {  
308         "words": [  
309           {  
310             "$word_context": "$word",  
311             "$begin_time": "$begintime",  
312             "end_time": "$endtime",  
313             "channel_id": "$channel",  
314             "word_index": "$index",  
315           }  
316         ]  
317       }  
318     }  
319   }  
320 }
```

318 E Datasheets for SpokenWOZ

319 E.1 Dataset documentation and intended uses

320 **For what purpose was the dataset created? Was there a specific task in mind? Was there a**
321 **specific gap that needed to be filled? Please provide a description.** Task-oriented dialogue (TOD)
322 models have made significant progress in recent years. These systems are designed to assist users in
323 accomplishing specific goals, e.g., flight booking and restaurant reservation. However, these TOD
324 datasets constructed solely based on written texts may not accurately reflect the nuances of spoken
325 conversations, leading to a gap between academic research and real-world spoken TOD scenarios.
326 We introduce the common tasks of TOD, including dialogue state tracking, policy Optimization, and
327 End-to-end Modeling.

328 **Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g.,**
329 **15 company, institution, organization)?** This dataset is created by researchers at Alibaba Group,
330 Renmin University of China, and University of Michigan.

331 **Who funded the creation of the dataset?** The creation of dataset was funded by DAMO Academy,
332 Alibaba Group.

333 **Any other comments?** N/A

334 E.2 Composition

335 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,**
336 **countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people**
337 **and interactions between them; nodes and edges)? Please provide a description.** The dataset
338 contains text data and audio data of spoken task-oriented dialogue. For each dialogue text, we also
339 annotate both the dialogue state and dialogue act. For the dialogue audio, we also give the ASR
340 transcription and audio file.

341 **How many instances are there in total (of each type, if appropriate)?** SpokenWOZ contains 8
342 domains, 203k turns, 5.7k dialogues and 249 hours of audios from spoken conversations.

343 **What data does each instance consist of? "Raw" data (e.g., unprocessed text or images) or fea-**
344 **tures? In either case, please provide a description.** For a conversation data, it includes a text part
345 and a audio part. For each dialogue text, SpokenWOZ contains dialogue context, annotated dialogue

346 state and annotated dialogue act. For the dialogue audio, SpokenWOZ contains corresponding audio
347 data and ASR transcription.

348 **Is there a label or target associated with each instance? If so, please provide a description** Yes,
349 we inherit and extend the MultiWOZ annotation schema, which is widely used for task-oriented
350 dialogue.

351 **Is any information missing from individual instances? If so, please provide a description,**
352 **explaining why this information is missing (e.g., because it was unavailable). This does not**
353 **include intentionally removed information, but might include, e.g., redacted text.** For individual
354 instances, there is no missing information.

355 **Are relationships between individual instances made explicit (e.g., users' movie ratings, social**
356 **network links)? If so, please describe how these relationships are made explicit.** Each dialogue
357 in SpokenWOZ is relatively independent, and the domains involved are different.

358 **Are there recommended data splits (e.g., training, development/validation, testing)? If so, please**
359 **provide a description of these splits, explaining the rationale behind them.** We give the data
360 splits in Appendix A.1.5. The data is split into training, development, and unreleased test sets. Once
361 researchers have built a model that works to your expectations on the dev set, they can submit it to us
362 to get official scores on the hidden test set. To mitigate the misestimation of the generalization error
363 of the model, we increase the number of test set to 1000 dialogues. At the same time we keep the
364 training and test data domain distributions roughly, but not exactly, the same.

365 **Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide**
366 **a description.** Since our dataset requires manual annotation of dialogue state and dialogue act,
367 annotation noise is inevitably introduced. At the same time the dialogue audio collection there
368 are cases of substandard audio quality, such as low communication quality. As shown in section
369 SpokenWOZ Construction, strict quality control is performed at each collection stage.

370 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**
371 **websites, tweets, other datasets)?** SpokenWOZ is self-contained.

372 **Does the dataset contain data that might be considered confidential (e.g., data that is pro-**
373 **ected by legal privilege or by doctor-patient confidentiality, data that includes the content of**
374 **individuals' non-public communications)? If so, please provide a description.** No.

375 **Does the dataset relate to people? If not, you may skip the remaining questions in this section.**
376 No. Detailed in Ethics Statement, we have designed the scenarios where an agent needs to proactively
377 collects user information, however, our user information is all generated by scripts in a andom manner,
378 so there will not be any privacy leakage issues.

379 E.3 Collection process

380 **How was the data associated with each instance acquired? Was the data directly observable**
381 **(e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly**
382 **inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or**
383 **language)? If data was reported by subjects or indirectly inferred/derived from other data, was**
384 **the data validated/verified? If so, please describe how.** We report the construction schema of
385 SpokenWOZ in the section SpokenWOZ Construction.

386 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or**
387 **sensor, manual human curation, software program, software API)? How were these mechanisms**
388 **or procedures validated?** We organized 250 participants to generate 5,700 dialogues via phone

389 calls. The details of the audio file can be found in Appendix A.1.4. The dialogue state and dialogue
390 act are annotated using the Appen platform².

391 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
392 **probabilistic with specific sampling probabilities)?** SpokenWOZ is not sampled from a larger
393 set.

394 **Over what timeframe was the data collected? Does this timeframe match the creation timeframe**
395 **of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please**
396 **describe the time-frame in which the data associated with the instances was created.** Our data
397 collection starts in July 2022. The contents of our data instances are independent of the time of
398 collection.

399 **Were any ethical review processes conducted (e.g., by an institutional review board)? If so,**
400 **please provide a description of these review processes, including the outcomes, as well as a**
401 **link or other access point to any supporting documentation.** Not applicable. We consider these
402 contents in Ethics Statement.

403 **Does the dataset relate to people? If not, you may skip the remainder of the questions in this**
404 **section.** No.

405 E.4 Preprocessing/cleaning/labeling

406 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,**
407 **tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing**
408 **of missing values)? If so, please provide a description. If not, you may skip the remainder of**
409 **the questions in this section.** We report the construction schema of SpokenWOZ in the section
410 SpokenWOZ Construction.

411 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support**
412 **unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.**
413 Raw data was not saved to prevent misuse. We will only open source the cleaned data.

414 **Is the software used to preprocess/clean/label the instances available? If so, please provide a**
415 **link or other access point.** We have used Python language to implement data cleaning. We will
416 share the scripts details in our codebase.

417 E.5 Uses

418 **Has the dataset been used for any tasks already? If so, please provide a description.** The
419 complexity and diverse spoken characteristics in SpokenWOZ make it a useful dataset for different
420 TOD tasks, including dialogue state tracking and response generation. For response generation, the
421 challenges are twofold: Policy Optimization and End-to-end Modeling. More details can be found in
422 section Tasks & Settings.

423 **Is there a repository that links to any or all papers or systems that use the dataset? If so, please**
424 **provide a link or other access point.** We provide links to the papers of all the baseline models on
425 the leaderboard³ we built.

426 **What (other) tasks could the dataset be used for?** The dataset can be used for the full range of
427 tasks related to task-oriented dialogue and can be used for dual-modal task-oriented dialogue studies.

²<https://appen.com/>

³<https://spokenwoz.github.io/SpokenWOZ-github.io/>

428 **Is there anything about the composition of the dataset or the way it was collected and prepro-**
429 **cessed/cleaned/labeled that might impact future uses? For example, is there anything that a**
430 **future user might need to know to avoid uses that could result in unfair treatment of individuals**
431 **or groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial**
432 **harms, legal risks) If so, please provide a description. Is there anything a future user could do**
433 **to mitigate these undesirable harms?** NA.

434 **Are there tasks for which the dataset should not be used? If so, please provide a description.**
435 NA.

436 E.6 Distribution

437 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institu-**
438 **tion, organization) on behalf of which the dataset was created?If so, please provide a de-**
439 **scription.** SpokenWOZ dataset and codebases for reproducing the experiments are available at:
440 <https://spokenwoz.github.io/SpokenWOZ-github.io/>.

441 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset**
442 **have a digital object identifier (DOI)?** The dataset is now available at: <https://spokenwoz.github.io/SpokenWOZ-github.io/>.

444 **When will the dataset be distributed?** The dataset is available now.

445 **Will the dataset be distributed under a copyright or other intellectual property (IP) license,**
446 **and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and**
447 **provide a link or other access point to, or otherwise reproduce, any relevant licensing terms**
448 **or ToU, as well as any fees associated with these restrictions.** SpokenWOZ is distributed under
449 the CC BY-NC 4.0⁴ license. CC BY-NC 4.0 allows reusers to distribute, remix, adapt, and build
450 upon the material in any medium or format for noncommercial purposes only, and only so long as
451 attribution is given to the creator. If you remix, adapt, or build upon the material, you must license
452 the modified material under identical terms.

453 **Have any third parties imposed IP-based or other restrictions on the data associated with the**
454 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
455 **or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these**
456 **restrictions.** No.

457 **Do any export controls or other regulatory restrictions apply to the dataset or to individual**
458 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
459 **or otherwise reproduce, any supporting documentation.** No.

460 E.7 Maintenance

461 **Who is supporting/hosting/maintaining the dataset?** Authors of this work bear all respon-
462 sibility in case of violation of rights. Shuzheng Si (sishuzheng@foxmail.com) and Wentao Ma
463 (mawentao.mwt@alibaba-inc.com) will be responsible for maintaining this dataset.

464 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?** If you
465 wish to extend or contribute to our dataset SpokenWOZ, please contact us via email - Shuzheng Si
466 (sishuzheng@foxmail.com) and Wentao Ma (mawentao.mwt@alibaba-inc.com).

467 **Is there an erratum? If so, please provide a link or other access point.** Any updates to the
468 dataset will be shared via GitHub

⁴<https://creativecommons.org/licenses/by-nc/4.0/legalcode>

469 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?**
470 **If so, please describe how often, by whom, and how updates will be communicated to users**
471 **(e.g., mailing list, GitHub)?** If we find inconsistencies in the dataset or extend the dataset, we will
472 release the new version on the website and Github.

473 **If the dataset relates to people, are there applicable limits on the retention of the data associated**
474 **with the instances (e.g., were individuals in question told that their data would be retained for a**
475 **fixed period of time and then deleted)?** N/A

476 **Will older versions of the dataset continue to be supported/hosted/maintained? If so, please**
477 **describe how. If not, please describe how its obsolescence will be communicated to users.** All
478 versions of SpokenWOZ will be continue to be supported and maintained on website. We will post
479 the updates on the website and Github.

480 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**
481 **them to do so? If so, please provide a description. Will these contributions be validated/verified?**
482 **If so, please describe how. If not, why not? Is there a process for communicating/distributing**
483 **these contributions to other users? If so, please provide a description.** Yes. Please contact the
484 authors of this paper for building upon this dataset.

485 E.8 Responsibility

486 The authors bear all responsibility in case of violation of rights, etc. We confirm that the dataset is
487 licensed under CC BY-NC 4.0 license.

488 References

- 489 [1] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy
490 Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. A
491 multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and
492 interactivity. CoRR, abs/2302.04023, 2023.
- 493 [2] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, M. Zhou,
494 and Hsiao-Wuen Hon. Unified language model pre-training for natural language understanding
495 and generation. ArXiv, abs/1905.03197, 2019.
- 496 [3] Vojtech Hudecek and Ondrej Dusek. Are llms all you need for task-oriented dialogue? CoRR,
497 abs/2304.06556, 2023.
- 498 [4] Benjamin Marie and Atsushi Fujita. Synthesizing parallel data of user-generated texts with zero-
499 shot neural machine translation. Transactions of the Association for Computational Linguistics,
500 8:710–725, 2020.
- 501 [5] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,
502 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton,
503 Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano,
504 Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback.
505 In NeurIPS, 2022.
- 506 [6] Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong
507 Chen. MultiWOZ 2.2 : A dialogue dataset with additional annotation corrections and state
508 tracking baselines. In Proceedings of the 2nd Workshop on Natural Language Processing for
509 Conversational AI, pages 109–117, Online, July 2020. Association for Computational Linguistics.