

DialogStudio: Towards Richest and Most Diverse Unified Dataset Collection for Conversational AI

Anonymous ACL submission

Abstract

Despite advancements in conversational AI, language models encounter challenges to handle diverse conversational tasks, and existing dialogue dataset collections often lack diversity and comprehensiveness. To tackle these issues, we introduce DialogStudio: the largest and most diverse collection of dialogue datasets, unified under a consistent format while preserving their original information. Our collection encompasses data from open-domain dialogues, task-oriented dialogues, natural language understanding, conversational recommendation, dialogue summarization, and knowledge-grounded dialogues, making it an incredibly rich and diverse resource for dialogue research and model training. To further enhance the utility of DialogStudio, we identify the licenses for each dataset, design external knowledge and domain-aware prompts for selected dialogues to facilitate instruction-aware fine-tuning. Furthermore, we develop conversational AI models using the dataset collection, and our experiments in both zero-shot and few-shot learning scenarios demonstrate the superiority of DialogStudio. To improve transparency and support dataset and task-based research, as well as language model pre-training, all datasets, licenses, codes, and models associated with DialogStudio will be made publicly accessible.¹

1 Introduction

Recent years have seen remarkable progress in Conversational AI, primarily driven by the advent of language models (Longpre et al., 2023; Zhang et al., 2022b; Brown et al., 2020; Touvron et al., 2023). Despite the advancements, these models could fall short when handling various tasks in a conversation due to the lack of comprehensive and diverse training data. Current dialogue

datasets (Lin et al., 2021; Asri et al., 2017) are typically limited in size and task-specific, which thus results in suboptimal ability in task-oriented model performance. Additionally, the lack of dataset standardization impedes model generalizability.

A few recent works (Gupta et al., 2022; Longpre et al., 2023; Ding et al., 2023) have introduced a large collection of datasets, which includes diverse tasks based on public datasets. For instance, FlanT5 (Longpre et al., 2023) presents the flan collections with a wide array of datasets and tasks. However, only a few of them are relevant to conversational AI. Although OPT (Iyer et al., 2022) have incorporated collections with several dialogue datasets, these collections remain inaccessible to the public. In contrast, efforts like InstructDial (Gupta et al., 2022) and ParLAI (Miller et al., 2017) consist of more dialogue datasets, but they lack diversity and comprehensiveness. For instance, ParLAI mainly includes open-domain dialogue datasets, which are exclusively accessible through their platform. Other collections (Gupta et al., 2022; Kim et al., 2022a; Ding et al., 2023; Dubois et al., 2023) often distill single dataset from ChatGPT or process datasets into a sequence-to-sequence format to support language model training, featuring input-output pairs such as dialogue context and system response. However, they often overlook other crucial dialogue information, constraining their utility for research interest on individual datasets, tasks, and broader applications.

To overcome the aforementioned challenges, we introduce DialogStudio, the most comprehensive and diverse collection of publicly available dialogue datasets, unified under a consistent format. By aggregating dialogue from various sources, DialogStudio promotes holistic analysis and the development of models adaptable to a variety of conversational scenarios. The collection spans an ex-

¹Due to the extensive size (~50GB) of our data and code, they will be made publicly available upon the paper's publication.

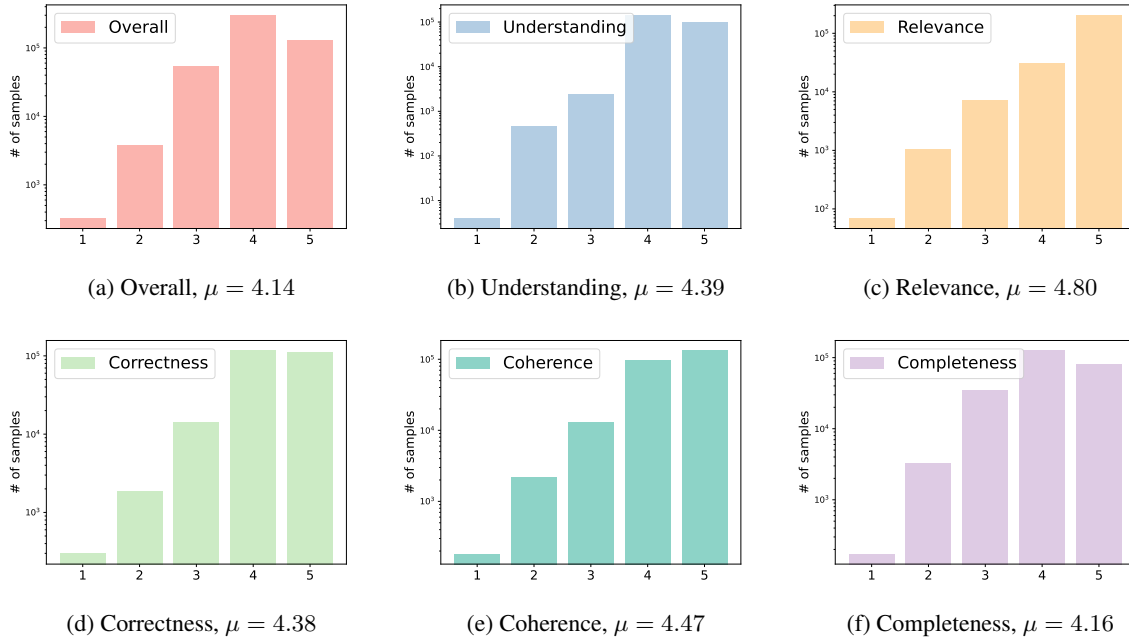


Figure 2: The score distribution for the dialogue quality.

138 *Summarization (Dial-Sum)*, and *Knowledge-*
 139 *Grounded Dialogues (KG-Dial)*. Figure 1a
 140 presents an overview of DialogStudio’s dataset
 141 categories. Note that the category boundaries are
 142 fuzzy as some datasets span multiple categories.
 143 For instance, SalesBot (Chiu et al., 2022) contains
 144 both casual and task-oriented conversations.
 145 Analogously, MultiWOZ (Budzianowski et al.,
 146 2018), a task-oriented dialogue corpus, incor-
 147 porates knowledge bases and dialogue acts to
 148 enhance knowledge-grounded generation. Addi-
 149 tionally, DialogStudio demonstrates its diversity
 150 by covering a wide range of domains, part of
 151 which is shown in Figure 1b.

152 2.2 Data Quality Investigation

153 Due to the existence of noise in dialogue, we
 154 develop a simple yet effective way to verify the
 155 quality of the datasets. Specifically, we employ
 156 ChatGPT (GPT-3.5-turbo) to evaluate the quality
 157 of system responses based on several perspec-
 158 tives (Mehri et al., 2022; Kim et al., 2022a), *i.e.*,
 159 Understanding, Relevance, Correctness, Coher-
 160 ence, Completeness and Overall quality. Under-
 161 standing assesses whether the model’s responses
 162 accurately reflect the meaning and intent of the
 163 user’s inputs. Relevance demonstrates whether the
 164 generated response should be directly related and
 165 appropriate to the preceding user input and the
 166 context of the conversation. Coherence measures

the logical consistency of the model’s responses
 within the context of the conversation. Comple-
 teness refers to whether the system’s responses fully
 address the user’s queries or tasks. Overall quality
 comprehensively rates the quality of dialogue. All
 scores are in the range of 1-5, and higher scores
 should only be given to truly exceptional exam-
 ples. We delicately design the prompt and ask the
 ChatGPT model to *strictly* rate the score.

Since there are a lot of datasets in DialogStu-
 dio, we randomly select 33 multi-turn dialogue
 datasets and evaluate all the training dialogues of
 each dataset. To harmonize ChatGPT and human
 ratings, we take a random sample of 50 training di-
 alogues from each dataset. These were then rated
 by three expert researchers using the five specified
 criteria. Post-alignment of ChatGPT and human
 evaluations, we view dialogues with a score above
 3 as being of high quality. Figure 2 illustrates dis-
 tributions of those scores. We also reveal the aver-
 age score as the μ in each sub-caption. In general,
 the dialogues show high qualities regarding to the
 individual criteria and the overall quality.

190 3 Datasets Unification and Access

191 We collect and process a wide range of datasets,
 192 involving different domains, types, and tasks.
 193 Since these datasets originally contain various in-
 194 formation and format, we propose a unification
 195 strategy to process all the datasets such that they

can be loaded in the same data loader.

3.1 Unification

Before unifying the format of those datasets, we fixed several issues as follows: 1) we remove those dialogues labeled as multi-turn dialogues, but actually with only one turn and miss either user utterance or system utterance. 2) We manually check the individual dialogues. If one dialogue contains one or more empty user or system utterances, we fill utterances based on corresponding dialogue contexts, dialogue acts, and dialogue information. In total, less than 0.5% of dialogues had these issues. To support research interest on individual datasets, we have flagged and rectified these problematic dialogues.

Additionally, we recognize the success of instruction tuning for dialogue models and thus we manually pre-define five different prompt templates for multi-turn dialogue datasets, such as *This is a bot helping users to {Task_Domain}. Given the dialogue context and external database, please generate a relevant system response for the user.* The *{Task_Domain}* is associated with the dialogue domain and we manually create a corresponding description. For example, if a dialogue is of domain *travel*, we set *{Task_Domain}* as *book a trip*. A concrete example of the prompt is demonstrated in Figure 3. Moreover, many datasets lack a direct mapping between dialogues and their domain information. To address this, we determine the domain of each dialogue using its intent, schema, APIs, and associated databases.

Next, we construct a uniform JSON dictionary format to store all relevant information of each dialogue as illustrated in Figure 3.² Compared with existing works, DialogStudio covers more dialogue information and is easier to retrieve the information for arbitrary dialogue-related tasks. Concretely, we include all dialogue-related information, such as the dialogue ID, data split label, domain, task, and content. Additionally, we identify the external knowledge, dialogue state tracking (DST) knowledge, and intent knowledge in the dialogue, which are the most beneficial knowledge for a dialogue.

Regarding external knowledge, we construct it based on information such as databases and dialogue acts. Since each dialogue dataset focuses

²More examples are available in the supplementary data materials.

on specific tasks or domains and has a different database and annotation schema, we unify such information into *external knowledge*. For example, if the user is looking for a hotel and asking for its address, the system response should be based on both the search results from the database and the dialogue context. To simulate the realistic situation and avoid directly providing the model with the ground truth resulting hotel, we also randomly sample four other candidate results and mix them with the ground truth result. All information is flattened and converted into a string as external knowledge.

To complete tasks and generate coherent responses, a dialogue system needs to track users' requirements for the task. Those requirements are usually represented as dialogue states. For example, regarding the hotel booking task, a dialogue system needs to extract information such as price range and locations to enable searching hotels in the database. The type of dialogue states varies across different tasks and datasets. As such, it is hard for dialogue systems to predict the values of those dialogue states if unknowing the specific dialogue states the task covers. Therefore, we propose to insert the schema, consisting of pre-defined dialogue state types and values for each task, into the input sequence. For datasets like SGD (Rastogi et al., 2020), which already provides annotation schema, we directly convert the dictionary-structured schema into a string. For the rest datasets that have no such schema file, we iterate over all dialogues and collect potential state annotations to construct a schema. We provide domains, slot types, and slot values in the schema string. For those categorized dialogue slots like "hotel star-ratings", which have a fixed number of candidate values, we provide all possible values. For others that have unlimited possible values, e.g. "stay night", we randomly sample ten values, such that a model can learn what slot values are relevant to these slot types. We put the turn-level ground-truth DST information in "dst", and the general DST information under "dst knowledge", as presented in Figure 3.

Analogously, intent prediction also requires models to know all possible intent types for each task. Therefore, we extract the schema directly from the schema file if it exists. As to datasets without schema, we also iterate over all dialogue in the dataset to collect all potential intents. Then,


```

"dialogue_id": "train_1",
"num_utterances": 14,
"utterances": [
  {
    "speaker": "USR",
    "text": "I'd like to book a trip to Atlantis from Caprica on
            Saturday, August 13, 2016 for 8 adults.",
    "ap_label": "",
    "da_label": "inform"
  },
  {
    "speaker": "USR",
    "text": "I have a tight budget of 1700.",
    "ap_label": "",
    "da_label": "inform"
  },
  {
    "speaker": "SYS",
    "text": "Hi...I checked a few options for you, and we do
            not currently have any trips that meet this criteria.",
    "ap_label": "",
    "da_label": "sorry",
    "slots": {
      "dst_city": "Atlantis",
      "or_city": "Caprica",
      "str_date": "Saturday, August 13, 2016",
      "n_adults": "8",
      "budget": "1700"
    }
  }
],
"scenario": {
  "db_id": "U22HTHYNP",
  "db_type": "booking",
  "task": "book"
}

```

(a) Original Data

```

"FRAMES--train--1": {
  "original_dialog_id": "train_1",
  "dialog_index": 1,
  "original_dialog_info": {
    "scenario": {
      "db_id": "U22HTHYNP",
      "db_type": "booking",
      "task": "book"
    }
  },
  "log": [
    {
      "turn_id": 1,
      "user_utterance": "I'd like to book a trip to Atlantis from Caprica on Saturday,
                        August 13, 2016 for 8 adults. I have a tight budget of 1700.",
      "system_response": "Hi...I checked a few options for you, and we do not currently
                          have any trips that meet this criteria.",
      "dialog_history": "",
      "original_user_side_information": {
        "da_label": "inform"
      },
      "original_system_side_information": {
        "da_label": "sorry",
        "slots": {
          "dst_city": "Atlantis",
          "or_city": "Caprica",
          "str_date": "Saturday, August 13, 2016",
          "n_adults": "8",
          "budget": "1700"
        }
      },
      "intent": "inform",
      "dst": "book dst_city Atlantis, book or_city Caprica, book str_date Saturday, August
             13, 2016, book n_adults 8, book budget 1700"
    }
  ],
  "external_knowledge": "( (travel : ((trip : (returning : (duration : (hours : 0 | min : 51...)",
  "dst_knowledge": "( (book : (dst_city : (Indianapolis | St. Louis | Le Paz | ...) | or_city : (
                    PUebLa | sf | toluca | San Francisco...)",
  "intent_knowledge": "( (book : (null | negate | request | goodbye | affirm))...)",
  "prompt": [
    "This is a bot helping users to book a trip. Given the dialog context and external
     database, please generate a relevant system response for the user."
  ]
}
}

```

(b) DialogStudio Data

Figure 3: A dialogue format example. Left: original example, right: converted example. Here we only show the first turn and partial information.

we put the turn-level ground-truth intent information into "intent", and the general intents under "intent knowledge", as presented in Figure 3. Note that not all datasets provide detailed annotation for dialogue states, intents, or even databases. For dialogue state tracking and intent classification tasks, we only process dialogues with corresponding annotations. Since all data is used for response generation, we leave the external knowledge value for the database blank if there is no related database in the original dataset.

3.2 Access and Maintenance

As aforementioned in the format, our DialogStudio data is easy to access via the JSON files. To make DialogStudio more maintainable and accessible, we will publish datasets on both GitHub and HuggingFace. GitHub mainly stores selected dialogue examples and relevant documents. We sample five original dialogues and five converted dialogues for each dataset to facilitate users in comprehending our format and examin-

ing the contents of each dataset. The complete DialogStudio dataset is maintained in our HuggingFace repository, where all the datasets can be directly downloaded or loaded with the HuggingFace `load_dataset(datasetname)` API. Given the substantial volume of datasets, optimizing user experience poses a challenge and limitation. We will continuously maintain and update both GitHub and HuggingFace.

DialogStudio is built upon public research datasets without individual or private information. We believe it is important to clearly present the license associated with each of these datasets. Consequently, we have included the original licenses for all datasets. All these datasets are supportive of academic research, and some of them also endorse commercial usage. The code that we employ falls under the widely accepted Apache 2.0 license. While we strictly require adherence to the respective dataset licenses for all intended usages on DialogStudio, there remains a possibility that some works might not fully comply with the li-

censes.

Regarding the other concerns such as ethical concern, we admit that DialogStudio is collected and maintained by the authors of this work and we did not hire external annotators. Since it contains unified datasets across several categories, it supports various research purposes from individual tasks and datasets to language model pre-training.

4 Experiments

In this section, we present the pre-training details, methodologies, and metrics used to assess the performance of our DialogStudio model. The evaluation process aims to measure the model’s ability to both solve task-oriented dialogues and understand general prompt-based instruction.

4.1 Model Pre-training

In this section, we introduce more details about how we conduct our pre-training. In regards of training models, we mix several datasets from DialogStudio.

For task-oriented and conversational recommendation datasets, we selected dialogues from a range of sources including KVRET (Eric et al., 2017), AirDialogue (Wei et al., 2018), DSTC2-Clean (Mrkšić et al., 2017), CaSiNo (Chawla et al., 2021), FRAMES (El Asri et al.), WOZ2.0 (Mrkšić et al., 2017), CraigslistBargains (He et al., 2018), Taskmaster1-2 (Byrne et al., 2019), ABCD (Chen et al., 2021a), MulDoGO (Peskov et al., 2019), BiTOD (Lin et al., 2021), SimJoint (Shah et al., 2018), STAR (Mosig et al., 2020), SGD (Rastogi et al., 2020), OpenDialog (Moon et al., 2019) and DuRecDial-2.0 (Liu et al., 2021).

Meanwhile, for knowledge-grounded dialogues, we drew upon dataset from SQA (Iyyer et al., 2017), SParC (Yu et al., 2019b), FeTaQA (Nan et al., 2022), MultiModalQA (Talmor et al., 2021), CompWebQ (Talmor and Berant, 2018), CoSQL (Yu et al., 2019a).

For open-domain dialogues, we sample dialogues from SODA (Kim et al., 2022a), Prosocial-Dialog (Kim et al., 2022b), Chitchat (Myers et al., 2020).

For each dialogue dataset, we sample at most 11k dialogues. Additionally, we extracted around 11k dialogue turns from question-answering dialogues featured in RACE (Lai et al., 2017), NarrativeQA (Kočiskỳ et al., 2018), SQUAD (Ra-

jpurkar et al., 2018), MCTest (Richardson et al., 2013), OpenBookQA (Mihaylov et al., 2018), MultiRC (Khashabi et al., 2018). Here, a dialogue turn refers to a pair consisting of a dialogue context and its corresponding system response. The rest datasets in DialogStudio are preserved for future evaluations and downstream fine-tuning.

For each dialogue during the training, we shape the available external knowledge into a string, which is included in dialogue context, and instruct the model to generate a dialogue response based on external knowledge. We use the format *Instruction* \n <USER> user utterance <SYSTEM> system response <USER> ... <USER> user utterance \n <EXTERNAL KNOWLEDGE> supported knowledge to train the model, where <USER>, <SYSTEM> and <EXTERNAL KNOWLEDGE> are special tokens.

We follow the public HuggingFace transformer code (Wolf et al., 2020; Wang et al., 2022) to train the model. For initializing our models, we adopt T5 (Raffel et al., 2020) and Flan-T5 (Longpre et al., 2023) as starting points to respectively establish DialogStudio-T5 and DialogStudio-Flan-T5. For the training of DialogStudio-Flan-T5, we exclude all translation-oriented tasks, limiting the sample size for each Flan task to a maximum of 150 examples. This leads to a cumulative total of 140,000 samples. We train the model up to 3 epochs with bfloat16 precision, a total batch size of 64. We set a constant learning rate 5e-5 and 3e-5 for the large model and the 3B model, respectively. Experiments are conducted using 16 A100 GPUs, each with 40GB of GPU memory.

4.2 Evaluation for Response Generation

Settings. We evaluate the performance on CoQA (Reddy et al., 2019) and MultiWOZ 2.2 (Zang et al., 2020). CoQA is a multi-turn conversational question answering dataset with 8k conversations about text passages from seven diverse domains. MultiWOZ 2.2 is one of the largest and most widely used multi-domain task-oriented dialogue corpora with more than 10000 dialogues. Each dialogue involves with one or more domains such as *Train, Restaurant, Hotel, Taxi, and Attraction*. The dataset is challenging and complex due to the multi-domain setting and diverse linguistic styles. Note that we exclude these two datasets during the pre-training stage in case of data leakage.

	CoQA		MultiWOZ	
	ROUGE-L	F1	ROUGE-L	F1
Flan-T5-3B (Longpre et al., 2023)	37.1	37.2	7.0	7.4
Flan-T5-Large (Longpre et al., 2023)	22.5	22.3	15.9	17.6
GODEL-Large (Peng et al., 2022)	43.2	43.3	18.5	19.3
DialogStudio-T5-Large	61.2	61.5	32.4	34.5
DialogStudio-Flan-T5-Large	63.3	63.5	33.7	35.9

Table 1: Zero-shot results on CoQA and MultiWOZ 2.2.

	CR (14 tasks)	DAR (7 tasks)	TE (27 tasks)	avg. (48 tasks)
OPT-30B (Zhang et al., 2022b)	21.3/1.1	35.2/4.1	40.3/0.9	32.3/2.0
OPT-IML-30B (Iyer et al., 2022)	37.4/41.6	51.4/51.8	54.7/47.8	47.9/47.1
OPT-175B (Zhang et al., 2022b)	21.0/4.2	37.1/16.8	41.6/2.2	33.3/7.7
OPT-IML-175B (Iyer et al., 2022)	39.0/49.8	61.2/60.2	54.3/51.0	51.5/53.6
Tk-INSTRUCT-11B (Wang et al., 2022)	32.3/ 62.3	51.1/ 69.6	55.0/64.1	46.1/ 65.3
Tk-INSTRUCT-3B (Wang et al., 2022)	38.4/51.3	45.7/58.5	48.4/52.8	44.2/54.2
DialogStudio-NIV2-T5-3B	41.3/59.8	57.5/63.7	52.3/55.1	50.4/59.5

Table 2: 0-shot/2-shot/5-shot ROUGE-L testing results on unseen datasets and unseen tasks. Results of baselines are reported by original papers. CR, DAR, and TE, avg. are abbreviations for Coreference Resolution, Dialogue Act Recognition, Textual Entailment, and average results, respectively.

For CoQA, we follow the original paper setting to answer question based on external passage. For MultiWOZ 2.2, we consider the lexicalized dialogue-act-to-response generation task where the model needs to consider both the dialogue context and the dialogue acts during generation. We follow the prompt from (Bang et al., 2023) to instruct models, i.e., *Continue the dialogue as a task-oriented dialogue system called SYSTEM. The answer of SYSTEM should follow the ACTION provided next while answering the USER’s last utterance.*

We focus on zero-shot evaluation and report the ROUGE-L and F1 score (Miller et al., 2017), where ROUGE-L measures the longest common subsequence and F1 measures the Unigram F1 overlap between the prediction and ground-truth response.

Baselines. We consider GODEL (Peng et al., 2022) and Flan-T5 (Longpre et al., 2023) as our baselines. GODEL is a T5-based large pre-trained model for goal-oriented dialogues. It is pre-trained with 551M multi-turn Reddit dialogues and 5M knowledge-grounded and question-answering dialogues. Flan-T5 is an instruction-aware model. It is also initialized from T5 and pre-trained on

the Flan collection, which consists of more than 1800 tasks and 400 datasets, including dialogue datasets.

Results. Table 1 depicts the results from both zero-shot and few-shot testing. Evidently, our models considerably surpass the baseline models in terms of zero-shot learning, exhibiting a robust generalized ability for response generation in a zero-shot scenario.

Flan-T5-3B, on the other hand, underperforms in the task of generating responses from dialog-acts. This model tends to produce incorrect dialog acts, unnatural utterances, or terminates with an empty end token. One explanation for this is that Flan-T5 models did not receive sufficient dialogue training during the instruction-training phase on the Flan collections. Comparisons between the performances of existing models before and after training on the unified dataset validate DialogStudio’s usefulness.

4.3 Evaluation on Super-NaturalInstructions

Settings. NIV2 (Wang et al., 2022) introduces an instruction-tuning benchmark with more than 1600 tasks. We select 3 categories with 44 tasks from the held-out test set, which consists of 154

	MMLU		BBH
	0-SHOT	5-SHOT	3-SHOT
TK-INSTRUCT 11B (Wang et al., 2022)	-	41.1	32.9
LLAMA 13B (Touvron et al., 2023)	-	46.2	37.1
Vicuna 13B (Chiang et al., 2023)	-	49.7	37.1
Flan-T5-Large (Longpre et al., 2023)	41.5	41.9	37.1
Flan-T5-XL (Peng et al., 2022)	48.7	49.3	40.2
OPT-IML-Max 30B (Iyer et al., 2022)	46.3	43.2	31.3
DialogStudio-Flan-T5-Large	40.1	40.9	34.2
DialogStudio-Flan-T5-3B	48.3	47.8	40.3

Table 3: Test results on MMLU and BBH. Results come from original papers and InstructEval (Chia et al., 2023).

tasks, i.e., Coreference Resolution, Dialogue Act Recognition, and Textual Entailment. The selected tasks and datasets are unseen in the training stage. Specifically, we follow all settings including metrics in (Wang et al., 2022), i.e., train models with instructions + 2 positive demonstrations and no negative demonstrations. We fine-tune DialogStudio-T5-3B on the 756 training tasks.

Baselines. OPT (Zhang et al., 2022b) is a set of open decoder-only transformer models pre-trained on 180B tokens. OPT-IML (Iyer et al., 2022) is an instruction meta-learning model based on the OPT-IML bench with more than 1500 tasks. Tk-INSTRUCT (Wang et al., 2022) is initialized from T5 and further pre-trained based on NIV2 collections. Note that we neglect Flan-T5 because it trains with all the downstream datasets and tasks.

Results. Table 2 shows the 0-shot and 2-shot results on unseen datasets and unseen tasks. Based on the average performance on 48 tasks, DialogStudio-NIV2-T5-3B outperforms OPT-IML-175B by 5.9% on 2-shot learning with more than 50 times fewer model parameters, and it surpasses Tk-INSTRUCT-11B by 4.3% on 0-shot learning with more than 3 times fewer parameters. The performance demonstrates the strong generalization ability of DialogStudio model. Compared with Tk-INSTRUCT-3B, DialogStudio-NIV2-T5-3B achieves 6.2% and 5.3% improvements on 0-shot and 2-shot learning respectively. Given that both Tk-INSTRUCT and our DialogStudio-NIV2-T5-3B are fine-tuned from the T5 model, this improvement indicates the effectiveness of pre-training with our DialogStudio collection.

4.4 Evaluation on MMLU and BBH

Table 3 presents results on MMLU (Hendrycks et al., 2020) and Big Bench Hard (BBH) (Srivastava et al., 2022).

When comparing the performance of DialogStudio-Flan-T5 with Flan-T5, we observe a minor decrease. However, this is accompanied by a significant improvement in dialogue relevant capabilities.

4.5 Evaluation on Alternative Benchmarks

DialogStudio encompasses not just public realistic dialogue datasets, but also those derived from or shared with ChatGPT, such as SODA (Kim et al., 2022a) and ShareGPT. Due to GPU constraints, we employ techniques like LoRA (Hu et al., 2021) to fine-tune llama (Touvron et al., 2023). When using equivalent datasets from DialogStudio, we observed performance comparable to other models, e.g., Vicuna (Chiang et al., 2023), on benchmarks like AlpacaEval (Dubois et al., 2023) and MT-Bench (Zheng et al., 2023). This demonstrates that DialogStudio caters to research interests in both specific datasets and generalized instruction tuning.

5 CONCLUSION

In this study, we have introduced DialogStudio, a comprehensive collection that aggregates more than 80 diverse dialogue datasets while preserving their original information. This aggregation not only represents a significant leap towards consolidating dialogues from varied sources but also offers a rich tapestry of conversational patterns, intents, and structures, capturing the nuances and richness of human interaction. Utilizing DialogStudio, we developed corresponding models, demonstrating superior performance in both zero-shot and few-shot learning scenarios. In the spirit of open research and advancing the field, we are committed to releasing DialogStudio to the broader research community.

561
562
563
564
565
566

567
568
569
570
571
572

573
574
575
576
577
578

579
580
581
582
583
584

585
586
587
588
589
590
591

592
593
594
595
596
597
598
599
600

601
602
603
604
605
606

607
608
609
610
611
612

613
614
615

References

Layla El Asri, Hannes Schulz, Shikhar Sharma, Jeremie Zumer, Justin Harris, Emery Fine, Rahul Mehrotra, and Kaheer Suleman. 2017. Frames: a corpus for adding memory to goal-oriented dialogue systems. *arXiv preprint arXiv:1704.00057*.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

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

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Ultes Stefan, Ramadan Osman, and Milica Gašić. 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 (EMNLP)*.

Bill Byrne, Karthik Krishnamoorthi, Chinnadhurai Sankar, Arvind Neelakantan, Ben Goodrich, Daniel Duckworth, Semih Yavuz, Amit Dubey, Kyu-Young Kim, and Andy Cediłnik. 2019. Taskmaster-1: Toward a realistic and diverse dialog dataset. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4516–4525.

Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. Efficient intent detection with dual sentence encoders. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45.

Iñigo Casanueva, Ivan Vulić, Georgios Spithourakis, and Paweł Budzianowski. 2022. Nlu++: A multi-label, slot-rich, generalisable dataset for natural language understanding in task-oriented dialogue. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1998–2013.

Kushal Chawla, Jaysa Ramirez, Rene Clever, Gale Lucas, Jonathan May, and Jonathan Gratch. 2021. Casino: A corpus of campsite negotiation dialogues

for automatic negotiation systems. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3167–3185.

Derek Chen, Howard Chen, Yi Yang, Alexander Lin, and Zhou Yu. 2021a. Action-based conversations dataset: A corpus for building more in-depth task-oriented dialogue systems. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3002–3017.

Maximillian Chen, Alexandros Papangelis, Chenyang Tao, Seokhwan Kim, Andy Rosenbaum, Yang Liu, Zhou Yu, and Dilek Hakkani-Tur. 2023. Places: Prompting language models for social conversation synthesis. In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 814–838.

Mingda Chen, Zewei Chu, Sam Wiseman, and Kevin Gimpel. 2022a. **SummScreen: A dataset for abstractive screenplay summarization**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8602–8615, Dublin, Ireland. Association for Computational Linguistics.

Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Yang Wang. 2020. Hybridqa: A dataset of multi-hop question answering over tabular and textual data. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1026–1036.

Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. 2021b. **DialogSum: A real-life scenario dialogue summarization dataset**. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 5062–5074, Online. Association for Computational Linguistics.

Zhiyu Chen, Bing Liu, Seungwan Moon, Chinnadhurai Sankar, Paul A Crook, and William Yang Wang. 2022b. Ketod: Knowledge-enriched task-oriented dialogue. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2581–2593.

Yew Ken Chia, Pengfei Hong, Lidong Bing, and Soujanya Poria. 2023. Instructeval: Towards holistic evaluation of instruction-tuned large language models. *arXiv preprint arXiv:2306.04757*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023).

Ssu Chiu, Maolin Li, Yen-Ting Lin, and Yun-Nung Chen. 2022. Salesbot: Transitioning from chit-chat to task-oriented dialogues. In *Proceedings of the*

672		Alexander Fabbri, Faiaz Rahman, Imad Rizvi, Borui Wang, Haoran Li, Yashar Mehdad, and Dragomir Radev. 2021. ConvoSumm: Conversation summarization benchmark and improved abstractive summarization with argument mining . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6866–6880, Online. Association for Computational Linguistics.	729
673			730
674			731
675	Samuel Coope, Tyler Farghly, Daniela Gerz, Ivan Vulić, and Matthew Henderson. 2020. Span-convert: Few-shot span extraction for dialog with pretrained conversational representations. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 107–121.		732
676			733
677			734
678			735
679			736
680			737
681	Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, et al. 2018. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. <i>arXiv preprint arXiv:1805.10190</i> .		738
682			
683			
684			
685			
686			
687			
688	Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander H. Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur D. Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhunoye, Alan W. Black, Alexander I. Rudnicky, Jason Williams, Joelle Pineau, Mikhail S. Burtsev, and Jason Weston. 2019. The second conversational intelligence challenge (convai2). <i>ArXiv</i> , abs/1902.00098.		739
689			740
690			741
691			742
692			743
693			744
694			745
695			746
696	Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2018. Wizard of wikipedia: Knowledge-powered conversational agents. <i>arXiv preprint arXiv:1811.01241</i> .		747
697			748
698			749
699			750
700	Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. <i>arXiv preprint arXiv:2305.14233</i> .		751
701			752
702			753
703			754
704			755
705	Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. AlpacaFarm: A simulation framework for methods that learn from human feedback. <i>arXiv preprint arXiv:2305.14387</i> .		756
706			757
707			758
708			759
709			760
710			761
711	Layla El Asri, Hannes Schulz, Shikhar Sharma, Jeremie Zumer, Justin Harris, Emery Fine, Rahul Mehrotra, and Kaheer Suleman. Frames: A corpus for adding memory to goal-oriented dialogue systems.		762
712			763
713			764
714			765
715			766
716	Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. Multiwoz 2.1: A consolidated multi-domain dialogue dataset with state corrections and state tracking baselines. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 422–428.		767
717			768
718			769
719			770
720			771
721			772
722			773
723			774
724	Mihail Eric, Lakshmi Krishnan, Francois Charette, and Christopher D Manning. 2017. Key-value retrieval networks for task-oriented dialogue. In <i>Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue</i> , pages 37–49.		775
725			776
726			777
727			778
728			779
			780
			781
			782
			783
			784
			785

786	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. <i>arXiv preprint arXiv:2009.03300</i> .	841
787		842
788		
789		
790	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	843
791		844
792		845
793		846
794		847
795		848
796	Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Dániel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, et al. 2022. Opt-impl: Scaling language model instruction meta learning through the lens of generalization. <i>arXiv preprint arXiv:2212.12017</i> .	849
797		850
798		851
799		852
800		853
801	Mohit Iyyer, Wen-tau Yih, and Ming-Wei Chang. 2017. Search-based neural structured learning for sequential question answering. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1821–1831.	854
802		855
803		856
804		857
805		858
806		859
807	Adam Janin, Don Baron, Jane Edwards, Dan Ellis, David Gelbart, Nelson Morgan, Barbara Peskin, Thilo Pfau, Elizabeth Shriberg, Andreas Stolcke, et al. 2003. The icsi meeting corpus. In <i>2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP'03)</i> , volume 1, pages I–I. IEEE.	860
808		861
809		862
810		863
811		864
812		865
813		866
814	Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: a challenge set for reading comprehension over multiple sentences. In <i>Proceedings of North American Chapter of the Association for Computational Linguistics (NAACL)</i> .	867
815		868
816		869
817		870
818		871
819		872
820	Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2022a. Soda: Million-scale dialogue distillation with social commonsense contextualization. <i>ArXiv</i> , abs/2212.10465.	873
821		874
822		875
823		876
824		877
825		878
826	Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. 2022b. Prosocialdialog: A prosocial backbone for conversational agents. <i>arXiv preprint arXiv:2205.12688</i> .	879
827		880
828		881
829		882
830		883
831	Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The narrativeqa reading comprehension challenge. <i>Transactions of the Association for Computational Linguistics</i> , 6:317–328.	884
832		885
833		886
834		887
835		888
836	Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2022. Internet-augmented dialogue generation. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8460–8478.	889
837		890
838		891
839		892
840		893
		894
		895
		896
	Wessel Kraaij, Thomas Hain, Mike Lincoln, and Wilfried Post. 2005. The ami meeting corpus.	
	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. Race: Large-scale reading comprehension dataset from examinations. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 785–794.	
	Stefan Larson, Anish Mahendran, Joseph J Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K Kummerfeld, Kevin Leach, Michael A Laurenzano, Lingjia Tang, et al. 2019. An evaluation dataset for intent classification and out-of-scope prediction. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1311–1316.	
	S Lee, H Schulz, A Atkinson, J Gao, K Suleman, L El Asri, M Adada, M Huang, S Sharma, W Tay, et al. 2019. Multi-domain task-completion dialog challenge. <i>Dialog system technology challenges</i> , 8(9).	
	Raymond Li, Samira Ebrahimi Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal. 2018a. Towards deep conversational recommendations. In <i>Advances in Neural Information Processing Systems 31 (NIPS 2018)</i> .	
	Xiujun Li, Yu Wang, Siqi Sun, Sarah Panda, Jingjing Liu, and Jianfeng Gao. 2018b. Microsoft dialogue challenge: Building end-to-end task-completion dialogue systems. <i>arXiv preprint arXiv:1807.11125</i> .	
	Yu Li, Kun Qian, Weiyang Shi, and Zhou Yu. 2020. End-to-end trainable non-collaborative dialog system. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 8293–8302.	
	Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, Peng Xu, Feijun Jiang, Yuxiang Hu, Chen Shi, and Pascale Fung. 2021. Bitod: A bilingual multi-domain dataset for task-oriented dialogue modeling. <i>NeurIPS 2021 Track on Datasets and Benchmarks</i> .	
	Jingjing Liu, Panupong Pasupat, Scott Cyphers, and Jim Glass. 2013. Asgard: A portable architecture for multilingual dialogue systems. In <i>2013 IEEE International Conference on Acoustics, Speech and Signal Processing</i> , pages 8386–8390. IEEE.	
	Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2019. Benchmarking natural language understanding services for building conversational agents. <i>arXiv preprint arXiv:1903.05566</i> .	
	Zeming Liu, Haifeng Wang, Zheng-Yu Niu, Hua Wu, and Wanxiang Che. 2021. Durecdial 2.0: A bilingual parallel corpus for conversational recommendation. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4335–4347.	

897	Shayne Longpre, Le Hou, Tu Vu, Albert Web-	Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Vic-	953
898	son, Hyung Won Chung, Yi Tay, Denny Zhou,	toria Lin, Neha Verma, Rui Zhang, Wojciech	954
899	Quoc V Le, Barret Zoph, Jason Wei, et al. 2023.	Kryściński, Hailey Schoelkopf, Riley Kong, Xian-	955
900	The flan collection: Designing data and methods	gru Tang, et al. 2022. Fetaqa: Free-form table ques-	956
901	for effective instruction tuning. <i>arXiv preprint</i>	tion answering. <i>Transactions of the Association for</i>	957
902	<i>arXiv:2301.13688</i> .	<i>Computational Linguistics</i> , 10:35–49.	958
903	Scott Martin, Shivani Poddar, and Kartikeya Upasani.	Linyong Nan, Dragomir Radev, Rui Zhang, Amrit	959
904	2020. Mudoco: corpus for multidomain coreference	Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xian-	960
905	resolution and referring expression generation. In	gru Tang, Aadit Vyas, Neha Verma, Pranav Krishna,	961
906	<i>Proceedings of the Twelfth Language Resources and</i>	et al. 2021. Dart: Open-domain structured data	962
907	<i>Evaluation Conference</i> , pages 104–111.	record to text generation. In <i>Proceedings of the 2021</i>	963
908	Shikib Mehri, Jinho Choi, Luis Fernando D’Haro, Jan	<i>Conference of the North American Chapter of the</i>	964
909	Deriu, Maxine Eskenazi, Milica Gasic, Kallirroi	<i>Association for Computational Linguistics: Human</i>	965
910	Georgila, Dilek Hakkani-Tur, Zekang Li, Verena	<i>Language Technologies</i> , pages 432–447.	966
911	Rieser, et al. 2022. Report from the nsf future direc-	Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann,	967
912	tions workshop on automatic evaluation of dialog:	Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and	968
913	Research directions and challenges. <i>arXiv preprint</i>	Dipanjan Das. 2020. Totto: A controlled table-to-	969
914	<i>arXiv:2203.10012</i> .	text generation dataset. In <i>Proceedings of the 2020</i>	970
915	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish	<i>Conference on Empirical Methods in Natural Lan-</i>	971
916	Sabharwal. 2018. Can a suit of armor conduct elec-	<i>guage Processing (EMNLP)</i> , pages 1173–1186.	972
917	tricity? a new dataset for open book question an-	Panupong Pasupat and Percy Liang. 2015. Compo-	973
918	swering. In <i>EMNLP</i> .	sitional semantic parsing on semi-structured tables.	974
919	Alexander H Miller, Will Feng, Adam Fisch, Jiasen Lu,	In <i>Proceedings of the 53rd Annual Meeting of the</i>	975
920	Dhruv Batra, Antoine Bordes, Devi Parikh, and Ja-	<i>Association for Computational Linguistics and the</i>	976
921	son Weston. 2017. Parlai: A dialog research soft-	<i>7th International Joint Conference on Natural Lan-</i>	977
922	ware platform. <i>arXiv preprint arXiv:1705.06476</i> .	<i>guage Processing (Volume 1: Long Papers)</i> , pages	978
923	Seungwhan Moon, Pararth Shah, Anuj Kumar, and Ra-	1470–1480.	979
924	jen Subba. 2019. Opendialkg: Explainable conversa-	Baolin Peng, Michel Galley, Pengcheng He, Chris	980
925	tional reasoning with attention-based walks over	Brockett, Lars Liden, Elnaz Nouri, Zhou Yu, Bill	981
926	knowledge graphs. In <i>Proceedings of the 57th An-</i>	Dolan, and Jianfeng Gao. 2022. Godel: Large-scale	982
927	<i>Annual Meeting of the Association for Computational</i>	pre-training for goal-directed dialog. <i>arXiv preprint</i>	983
928	<i>Linguistics</i> , pages 845–854.	<i>arXiv:2206.11309</i> .	984
929	Johannes EM Mosig, Shikib Mehri, and Thomas	Denis Peskov, Nancy Clarke, Jason Krone, Brigi Fodor,	985
930	Kober. 2020. Star: A schema-guided dialog	Yi Zhang, Adel Youssef, and Mona Diab. 2019.	986
931	dataset for transfer learning. <i>arXiv preprint</i>	Multi-domain goal-oriented dialogues (multidogo):	987
932	<i>arXiv:2010.11853</i> .	Strategies toward curating and annotating large scale	988
933	Nikola Mrkšić, Diarmuid Ó Séaghdha, Tsung-Hsien	dialogue data. In <i>Proceedings of the 2019 Confer-</i>	989
934	Wen, Blaise Thomson, and Steve Young. 2017.	<i>ence on Empirical Methods in Natural Language</i>	990
935	Neural belief tracker: Data-driven dialogue state	<i>Processing and the 9th International Joint Confer-</i>	991
936	tracking. In <i>Proceedings of the 55th Annual Meet-</i>	<i>ence on Natural Language Processing (EMNLP-</i>	992
937	<i>ing of the Association for Computational Linguistics</i>	<i>IJCNLP)</i> , pages 4526–4536.	993
938	(<i>Volume 1: Long Papers</i>), pages 1777–1788.	Kun Qian, Satwik Kottur, Ahmad Beirami, Shahin	994
939	Rajdeep Mukherjee, Abhinav Bohra, Akash Banerjee,	Shayandeh, Paul A Crook, Alborz Geramifard, Zhou	995
940	Soumya Sharma, Manjunath Hegde, Afreen Shaikh,	Yu, and Chinnadhurai Sankar. 2022. Database	996
941	Shivani Shrivastava, Koustuv Dasgupta, Niloy Gan-	search results disambiguation for task-oriented di-	997
942	guly, Saptarshi Ghosh, and Pawan Goyal. 2022.	alog systems. In <i>Proceedings of the 2022 Confer-</i>	998
943	ECTSum: A new benchmark dataset for bullet point	<i>ence of the North American Chapter of the Associ-</i>	999
944	summarization of long earnings call transcripts. In	<i>ation for Computational Linguistics: Human Lan-</i>	1000
945	<i>Proceedings of the 2022 Conference on Empirical</i>	<i>guage Technologies</i> , pages 1158–1173.	1001
946	<i>Methods in Natural Language Processing</i> , pages	Jun Quan, Deyi Xiong, Bonnie Webber, and Changjian	1002
947	10893–10906, Abu Dhabi, United Arab Emirates.	Hu. 2019. Gecor: An end-to-end generative ellipsis	1003
948	Association for Computational Linguistics.	and co-reference resolution model for task-oriented	1004
949	Will Myers, Tyler Etchart, and Nancy Fulda. 2020.	dialogue. In <i>Proceedings of the 2019 Conference on</i>	1005
950	Conversational scaffolding: An analogy-based ap-	<i>Empirical Methods in Natural Language Processing</i>	1006
951	proach to response prioritization in open-domain di-	<i>and the 9th International Joint Conference on Natu-</i>	1007
952	alogs.	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	1008
		4547–4557.	1009

1010	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>The Journal of Machine Learning Research</i> , 21(1):5485–5551.	<i>American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 641–651.	1066 1067 1068
1016	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for squad. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789.	Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hananeh Hajishirzi, and Jonathan Berant. 2021. Multimodalqa: Complex question answering over text, tables and images. <i>ICLR</i> .	1069 1070 1071 1072 1073
1021	Revanth Rameshkumar and Peter Bailey. 2020. Storytelling with dialogue: A Critical Role Dungeons and Dragons Dataset . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5121–5134, Online. Association for Computational Linguistics.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	1074 1075 1076 1077 1078 1079
1027	Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: a new benchmark and dataset. In <i>ACL</i> .	Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. 2022. Supernaturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5085–5109.	1080 1081 1082 1083 1084 1085 1086 1087
1031	Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 8689–8696.	Wei Wei, Quoc Le, Andrew Dai, and Jia Li. 2018. Airdialogue: An environment for goal-oriented dialogue research. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3844–3854.	1088 1089 1090 1091 1092
1037	Siva Reddy, Danqi Chen, and Christopher D Manning. 2019. Coqa: A conversational question answering challenge. <i>Transactions of the Association for Computational Linguistics</i> , 7:249–266.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations</i> , pages 38–45.	1093 1094 1095 1096 1097 1098 1099 1100
1041	Virgile Rennard, Guokan Shang, Julie Hunter, and Michalis Vazirgiannis. 2023. Abstractive meeting summarization: A survey. <i>Transactions of the Association for Computational Linguistics</i> , 11:861–884.	Wen-tau Yih, Matthew Richardson, Christopher Meek, Ming-Wei Chang, and Jina Suh. 2016. The value of semantic parse labeling for knowledge base question answering. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 201–206.	1101 1102 1103 1104 1105 1106
1045	Matthew Richardson, Christopher JC Burges, and Erin Renshaw. 2013. Mctest: A challenge dataset for the open-domain machine comprehension of text. In <i>Proceedings of the 2013 conference on empirical methods in natural language processing</i> , pages 193–203.	Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, et al. 2019a. Cosql: A conversational text-to-sql challenge towards cross-domain natural language interfaces to databases. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1962–1979.	1107 1108 1109 1110 1111 1112 1113 1114 1115
1051	Pararth Shah, Dilek Hakkani-Tür, Gokhan Tür, Abhinav Rastogi, Ankur Bapna, Neha Nayak, and Larry Heck. 2018. Building a conversational agent overnight with dialogue self-play. <i>arXiv preprint arXiv:1801.04871</i> .	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. 2018. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3911–3921.	1116 1117 1118 1119 1120 1121 1122 1123
1056	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>arXiv preprint arXiv:2206.04615</i> .		
1063	Alon Talmor and Jonathan Berant. 2018. The web as a knowledge-base for answering complex questions. In <i>Proceedings of the 2018 Conference of the North</i>		

1124 Tao Yu, Rui Zhang, Michihiro Yasunaga, Yi Chern
1125 Tan, Xi Victoria Lin, Suyi Li, Heyang Er, Irene Li,
1126 Bo Pang, Tao Chen, et al. 2019b. Sparc: Cross-
1127 domain semantic parsing in context. *arXiv preprint*
1128 *arXiv:1906.02285*.

1129 Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara,
1130 Raghav Gupta, Jianguo Zhang, and Jindong Chen.
1131 2020. Multiwoz 2.2: A dialogue dataset with addi-
1132 tional annotation corrections and state tracking base-
1133 lines. In *Proceedings of the 2nd Workshop on Nat-
1134 ural Language Processing for Conversational AI*,
1135 pages 109–117.

1136 Jianguo Zhang, Kazuma Hashimoto, Yao Wan, Zhiwei
1137 Liu, Ye Liu, Caiming Xiong, and S Yu Philip. 2022a.
1138 Are pre-trained transformers robust in intent classi-
1139 fication? a missing ingredient in evaluation of out-
1140 of-scope intent detection. In *Proceedings of the 4th*
1141 *Workshop on NLP for Conversational AI*, pages 12–
1142 20.

1143 Susan Zhang, Stephen Roller, Naman Goyal, Mikel
1144 Artetxe, Moya Chen, Shuohui Chen, Christopher
1145 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al.
1146 2022b. Opt: Open pre-trained transformer language
1147 models. *arXiv preprint arXiv:2205.01068*.

1148 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
1149 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
1150 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023.
1151 Judging llm-as-a-judge with mt-bench and chatbot
1152 arena. *arXiv preprint arXiv:2306.05685*.

1153 Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia
1154 Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli
1155 Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir
1156 Radev. 2021. QMSum: A New Benchmark for
1157 Query-based Multi-domain Meeting Summariza-
1158 tion. In *North American Association for Computa-
1159 tional Linguistics (NAACL)*.

1160 Victor Zhong, Caiming Xiong, and Richard Socher.
1161 2017. Seq2sql: Generating structured queries
1162 from natural language using reinforcement learning.
1163 *arXiv preprint arXiv:1709.00103*.

1164 Chenguang Zhu, Yang Liu, Jie Mei, and Michael Zeng.
1165 2021. Mediasum: A large-scale media interview
1166 dataset for dialogue summarization. *arXiv preprint*
1167 *arXiv:2103.06410*.

Appendix

Table 4 and Table 5 lists datasets included in DialogStudio. Initially, we present a partial list of these datasets. More dialogue examples are available in the supplementary data materials.

NLU	NLU++ (Casanueva et al., 2022)
	BANKING77-OOS (Zhang et al., 2022a)
	BANKING77 (Casanueva et al., 2020)
	RESTAURANTS8K (Coope et al., 2020)
	CLINC150 (Larson et al., 2019)
	CLINC-Single-Domain-OOS-banking (Zhang et al., 2022a)
	CLINC-Single-Domain-OOS-credit_cards (Zhang et al., 2022a)
	HWU64 (Liu et al., 2019)
	SNIPS (Coucke et al., 2018)
	SNIPS-NER (Coucke et al., 2018)
	DSTC8-SGD (Coope et al., 2020)
	TOP (Gupta et al., 2018)
	TOP-NER (Gupta et al., 2018)
	ATIS-NER (Hemphill et al., 1990)
	ATIS (Hemphill et al., 1990)
	MIT-MOVIE (Liu et al., 2013)
MIT-RESTAURANT (Liu et al., 2013)	
TOD	KVRET (Eric et al., 2017)
	AirDialogue (Wei et al., 2018)
	DSTC2-Clean (Mrkšić et al., 2017)
	CaSiNo (Chawla et al., 2021)
	FRAMES (El Asri et al.)
	WOZ2.0 (Mrkšić et al., 2017)
	CraigslistBargains (He et al., 2018)
	Taskmaster1 (Byrne et al., 2019)
	Taskmaster2 (Byrne et al., 2019)
	Taskmaster3 (Byrne et al., 2019)
	ABCD (Chen et al., 2021a)
	MulDoGO (Peskov et al., 2019)
	BiTOD (Lin et al., 2021)
	SimJointGEN (Shah et al., 2018)
	SimJointMovie (Shah et al., 2018)
	SimJointRestaurant (Shah et al., 2018)
	STAR (Mosig et al., 2020)
	SGD (Rastogi et al., 2020)
	MultiWOZ2.1 (Eric et al., 2020)
	MultiWOZ2.2 (Zang et al., 2020)
	HDSA-Dialog (Chen et al., 2021a)
	MS-DC (Li et al., 2018b)
	GECOR (Quan et al., 2019)
	Disambiguation (Qian et al., 2022)
	MetaLWOZ (Lee et al., 2019)
	KETOD (Chen et al., 2022b)
	MuDoCo (Martin et al., 2020)

Table 4: List of datasets included in DialogStudio (a).

KG-Dial	<p>SQA (Iyyer et al., 2017) SParC (Yu et al., 2019b) FeTaQA (Nan et al., 2022) MultiModalQA (Talmor et al., 2021) CompWebQ (Talmor and Berant, 2018) CoSQL (Yu et al., 2019a) CoQA (Reddy et al., 2019) Spider (Yu et al., 2018) ToTTo (Parikh et al., 2020) WebQSP (Yih et al., 2016) WikiSQL (Zhong et al., 2017) WikiTQ (Pasupat and Liang, 2015) DART (Nan et al., 2021) GrailQA (Gu et al., 2021) HybridQA (Chen et al., 2020) MTOPI (Chen et al., 2020) UltralChat-Assistance (Ding et al., 2023) Wizard_of_Wikipedia (Dinan et al., 2018) Wizard_of_Internet (Komeili et al., 2022)</p>
Dial-Sum	<p>TweetSumm (Feigenblat et al., 2021) SAMSum (Gliwa et al., 2019) DialogSum (Chen et al., 2021b) AMI (Kraaij et al., 2005; Rennard et al., 2023) ICSI (Janin et al., 2003) QMSum (Zhong et al., 2021) MediaSum (Zhu et al., 2021) ECTSum (Mukherjee et al., 2022) SummScreen.ForeverDreaming (Chen et al., 2022a) SummScreen.TVMegaSite (Chen et al., 2022a) CRD3 (Rameshkumar and Bailey, 2020) ConvoSumm (Fabbri et al., 2021)</p>
Open-Domain	<p>ChitCHAT (Myers et al., 2020) SODA (Kim et al., 2022a) Prosocial (Kim et al., 2022b) HH-RLHF (Bai et al., 2022) Empathetic (Rashkin et al., 2019) ConvAI2 (Dinan et al., 2019) AntiScam (Li et al., 2020) ShareGPT (Zheng et al., 2023) PLACES3.5 (Chen et al., 2023)</p>
Conv-Rec	<p>SalesBot (Chiu et al., 2022) Redial (Li et al., 2018a) Inspired (Hayati et al., 2020) DuRecDial 2.0 (Liu et al., 2021) OpendialKG (Moon et al., 2019)</p>

Table 5: List of datasets included in DialogStudio (b).