

UNO-DST: Leveraging Unlabelled Data in Zero-Shot Dialogue State Tracking

Anonymous ACL submission

Abstract

Previous zero-shot dialogue state tracking (DST) methods only apply transfer learning, but ignore unlabelled data in the target domain. We transform zero-shot DST into few-shot DST by utilising such unlabelled data via joint and self-training methods. Our method incorporates auxiliary tasks that generate slot types as inverse prompts for main tasks, creating slot values during joint training. Cycle consistency between these two tasks enables the generation and selection of quality samples in unknown target domains for subsequent fine-tuning. This approach also facilitates automatic label creation, thereby optimizing the training and fine-tuning of DST models. We demonstrate this method’s effectiveness on large language models in zero-shot scenarios, improving average joint goal accuracy by 8% across all domains in MultiWOZ¹

1 Introduction

Dialogue state tracking (DST) is a crucial task in understanding users’ intentions by extracting the dialogue states from the dialogue history (Balaraman et al., 2021), where a single dialogue state is a combination of a slot type (e.g., `<hotel-name>`) and a slot value (e.g., `<Hilton hotel>`), as in Figure 1. Dialogue states are a set of those combinations (e.g., `<hotel-name: Hilton hotel>`) retrieved by DST models, given dialogue history and slot types. Traditional methods train and evaluate DST models with manually-labelled dialogue states in each domain, which can be costly and time consuming (Wu et al., 2020b; Hosseini-Asl et al., 2020). Recently, DST under zero and few-shot settings draw increased attention (Lin et al., 2021b; Hudeček et al., 2021). Compared with few-shot methods, zero-shot approaches are more challenging, due to unseen slot types and data scarcity in unknown target domains.

¹Code and data are available at <https://anonymous.4open.science/r/UNO-DST-58B4>.

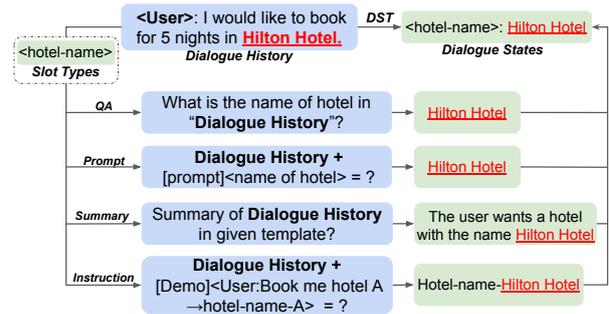


Figure 1: Examples of zero-shot methods in DST.

In both zero and few-shot settings, the majority of existing methods convert the DST problem into other common problem settings in natural language processing (NLP): for example, Question Answering (QA; Lin et al., 2021a; Li et al., 2021), prompt learning (Lee et al., 2021), summarization (Shin et al., 2022) and instruction learning (Gupta et al., 2022). For example in Figure 1, a given slot type (`<hotel-name>`) can be transformed into a QA setting by queries like “What is the hotel name in the context?” and the slot values can be predicted by a QA model accordingly.

Transfer learning methods also convert DST tasks to generation ones, more suited for pre-trained language models (PLMs; Devlin et al., 2019). However, such methods cannot fully leverage the capability of PLMs in generation and selection. Two main difficulties emerge: 1) the performance of the chosen NLP tasks can be unpredictable for unseen slot types in a new domain due to domain divergence; and 2) existing models are only trained in the known domains, without utilizing any unlabeled data in the new target domain.

This work proposes UNO-DST², a method to leverage the **unlabelled** data for **zero-shot DST** in target domain. Inspired by the popularity of multi-

²“Uno”, Spanish for “one”, embodies our proposed strategy in this paper: transitioning from zero to one and subsequently from one to all.

task learning and self-supervised learning (Zhang and Yang, 2021; Tsai et al., 2021), UNO-DST employs a two-step training framework invoking both joint and self-training (Figure 2). Aside from the main task of generating slot values, we design an auxiliary task of generating slot types. We then jointly train both tasks using the labelled training data in known source domains. For the self-training period, we implement the concept of cycle consistency within our two tasks (Zhu et al., 2017). That is, a text output from the main task serves as input to the auxiliary task, and the resultant text produced by the auxiliary task should match the original input text. This process forms a full cycle, ensuring consistent generation and selection of dialogue states from the unlabelled data, which will be further used for fine-tuning the model. In this way, we convert zero-shot problems into few-shot ones. Importantly, our framework is model-agnostic and can be applied to different baseline models.

Our main contributions are as follows:

- To the best of our knowledge, we are the first zero-shot DST work to use unlabelled training data in an unknown target domain under a two-step training strategy;
- We introduce an auxiliary task to facilitate the training of the main task, the selection of fine-tuning samples, and the generation of unseen or new slot types;
- Our results achieve a new state-of-the-art. Our analysis also identifies the lower and upper bounds of each zero-shot DST method.
- We demonstrate our methods with PLMs and large language models, showing its effectiveness on two popular DST datasets.

2 Related Works

Existing DST methods are generally classified as either 1) full-data or 2) low-resource DST. However, regardless of which method is adopted, unlabeled training data in the target domain is unutilized.

Full-data DST are commonly trained with fully annotated multi-domain conversations (Wu et al., 2020a; Hosseini-Asl et al., 2020). SOTA models focus on DST tasks with well-annotated datasets (Mrkšić et al., 2017; Ren et al., 2018). However, the annotation work for data in a new domain can be costly. Hence there is interest in transferring the knowledge of a model from a known domain into

an unknown domain and conducting DST tasks in a low-resource setting.

Low-resource DST uses zero- or few-shot learning in the unknown target domain. Here, the state-of-the-art use a single NLP task to transfer knowledge from the source domains to the unknown target domain (Lin et al., 2021b; Shin et al., 2022). While transfer learning tasks achieve good results, each method is task-dependent. Thus, task-independent strategies have been proposed (Wang et al., 2022; Yang et al., 2023).

Multi-task Learning involving simultaneous training of a model on diverse tasks, to boost performance on trained downstream tasks. It also holds promise for enhancements on new tasks (Raffel et al., 2020; Zhang and Yang, 2021). However, existing methods typically neglect to assess the consistency across multiple tasks after joint training, while our approach leverages the cycle consistency for selection (Zhu et al., 2017; Wang et al., 2023).

3 Methodology

In Figure 2, we show an overview of **UNO-DST** with joint training and self-training periods. Our method includes two tasks: a main task for slot value prediction (§3.1) and an auxiliary task for slot type prediction (§3.2). In the joint training period, both tasks are jointly trained in the known source domains (§3.3). In the self-training period, we introduce three steps to generate dialogue states, select good samples, and fine-tune the PLM (§3.4). Lastly, we elaborate on the transferability of our strategy with a lower and upper bound (§3.5).

3.1 Task Definition

The main task for DST is predicting the $\langle slot\text{-}type:slot\text{-}value \rangle$ pairs with given dialogue history and slot types from a pre-defined slot type list, as shown in Figure 1. For each domain, there are seen slot types which appear in other domains (e.g., “hotel-name” and “restaurant-name”) or unseen slot types which are unique in the specific domain (e.g., “hotel-stars”). The ratio of the occurrences of these unseen slot types represents the difficulty of zero-shot DST for each domain (Wang et al., 2022).

We denote the dialogue history in a t -turn conversations as $C_t = \{c_1, c_2 \dots c_t\}$ and slot types S in domain h as $S_h = \{s_1, s_2 \dots s_n\}$. For each conversation turn, the main goal is to predict slot values v' . Therefore, the input for the *PLM* is combined

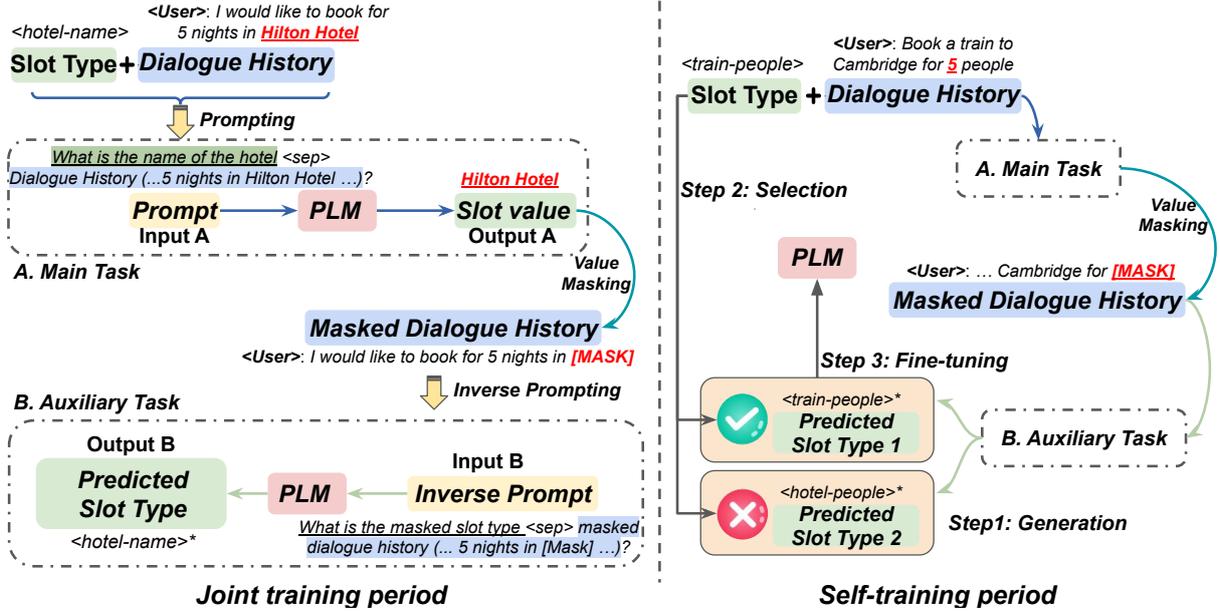


Figure 2: Overview of UNO-DST which consists of two periods: 1) joint training for both task A (slot value prediction) and B (slot type prediction), 2) self-training in the unknown target domain. Step 1: Generation of slot values and types from task A and B; Step 2: Selection of good samples using cycle consistency between two tasks; Step 3: Fine-tuning the PLM with selected samples.

of dialogue history and slot types, with the output being slot values, as shown in Eq. 1.

$$v'_i = PLM(C_t, s_i) \quad (1)$$

Compared with methods that select slot values from a constant ontology list using classification models (Shi et al., 2017), we enhance the capability of text-to-text PLMs for text generation (Heck et al., 2020). For the case when there are no slot values related to a given slot type, we train the model to output a “none” value, indicating that there are no dialogue states from the current conversation turn.

To better utilise the capability of PLMs in different tasks, we utilize different prompt functions “ $P(\cdot)$ ” to generate the prompt in correct format. For example, given a slot s and context c , the QA prompt p for the DST can be $p^{main} = \text{“What is the value of slot } s \text{ in context } c\text{?”}$. We formulate the way of using prompts for the DST main task as:

$$v'_i = PLM(p_i^{main}) = PLM(P(s_i, C_t)) \quad (2)$$

3.2 Auxiliary task

As joint training can improve the accuracy of PLMs, we design an auxiliary task to facilitate the training of the main task (Zhang and Yang, 2021; Su et al., 2022; Yang et al., 2023). We propose an auxiliary task to help the model better understand

the semantic and context information from the dialogue history in the joint training period and serve as a regulator to check the main task predictions obtained during the self-training period.

We design the auxiliary task as the inverse (converse) prompt of the main task. In opposition to the main task, the auxiliary task thus takes the slot values v as input and generates the slot types s' as outputs, which forms a cycle-consistent loop as a foil to the main task. To make it easier for PLMs, we convert the slot values v and dialogue history C_t into a masked dialogue history C_t^m for the model to make better masked predictions, as in Eq. 3. The inverse QA prompt p^{aux} is generated as “What is the masked slot type in context C_t^m ?” from inverse prompt function “ $IP(\cdot)$ ” (Figure 2). We implement the auxiliary task during both the joint training and self-training periods to facilitate slot values generation and selection.

$$s'_i = PLM(p_i^{aux}) = PLM(IP(v_i, C_t^m)) \quad (3)$$

3.3 Joint training with auxiliary tasks

We conduct a simple version of joint training with only two tasks: the main DST task and the auxiliary task. The training samples for the main tasks are created using a prompt of dialogue history and slot type, while the samples for the auxiliary DST tasks are created using an inverse prompt of masked dialogue history.

As the auxiliary task is an inverse of the main task, the model is trained for the same knowledge in a cycle-consistent way. By predicting the masked slot type from the masked dialogue history, the model is familiar with the context and different slot types. With our specially-designed auxiliary task, the generation model reuses the existing data for another round of training without the need to increase the amount of training data or model parameters. We formulate the loss function for the main task L_m and auxiliary task L_a as:

$$L_m = - \sum_i^n \log p(v'_i | C_t, s_i) \quad (4)$$

$$L_a = - \sum_i^n \log p(s'_i | C_t^m, v_i) \quad (5)$$

The final loss is a simple average of both. To keep the process simple, we do not add hyperparameters to the existing model framework. Importantly, as the auxiliary task samples are generated using the main task’s inverse prompt, the task ratio mirrors the natural distribution of both tasks throughout the joint training period.

3.4 Self-training with auxiliary tasks

Compared with other zero-shot DST models, the key novelty of our strategy is in using the unlabelled training data in the unknown target domain for self-training. Self-training aims to generate pseudo labels and select data samples that further fine-tune the models. In the self-training period, we divide the strategy into three steps: termed generation, selection and fine-tuning (Algorithm 1).

Step 1 Generation tests both tasks using the unlabelled training data in the unknown target domain to generate predicted slot values v' and slot types s' . Auxiliary tasks in self-training are created by value masking, as shown in Figure 2. For training samples with slot values that do not directly copy from the original context (such as “yes/no” for “hotel-parking”), masking the slot value in the original context does not work. Such samples are omitted in creating the masked dialogue history in auxiliary tasks.

Step 2 Selection tests the cycle consistency between main tasks and auxiliary tasks by comparing the predicted slot types s' with the original slot types s in each dialogue turn. A simplified selection process is shown in Figure 2. In experiments, only conversation turns with all correct slot types

are selected as good samples, similar to joint goal accuracy settings, aiming to reduce the selection error rate.

Step 3 fine-tunes the model $PLM(\cdot)$ with selected samples and predicted slot values v' . This completes the conversion of zero-shot DST into few-shot DST, helping the model adapt to unknown domains without increasing data annotation and model parameters. For models that are difficult to fine-tune, we propose other solutions (§ 7).

3.5 Lower and upper bound for zero-shot DST models

Even though there are many studies working on zero-shot DST, to the best of our knowledge there are no common methods to identify the threshold below where results are unreasonably low (lower bound) or the peak performance that each model can potentially achieve (upper bound). Here, we discuss our proposed algorithm with respect to the lower and upper bounds, aiming to limit the research that is unsuitable for zero-shot DST by the lower bound and benchmark our method against oracular results for each model as an upper bound.

Lower bound. PLMs generate either a “slot value” representing specific information, or “none” if no information is relevant to the dialogue history, for all slot types in our main tasks. Notably, most dialogues do not associate slot types with specific information, resulting in “none” being the prevalent prediction. For example, slot type “hotel-area” is linked to “none” for dialogue in Figure 2. Therefore, we set the lower bound as the outcome when models predict “none” for all slot types across all dialogue turns.

Upper bound. According to our self-training methods (§3.4), zero-shot DST can always be converted into a few-shot DST by selecting good samples with self-generated slot values for fine-tuning, but cycle consistency cannot ensure 100% correct data selection. Oracular performance comes when we select only the correct samples from all the self-generated slot values and use them to fine-tune the model. We define such performance as the upper bound for the zero-shot DST model.

4 Experiments and Datasets

Dataset. We train and test our model on both MultiWOZ 2.1 (Budzianowski et al., 2018) and the Schema-Guided Dialogue (SGD; Rastogi et al., 2020). MultiWOZ and SGD have dialogues dis-

	MultiWOZ		SGD	
	JGA		JGA	AGA
Benchmarks	25.8		27.6	58.0
T5DST	32.4		NA	NA
SD-T5	35.6		NA	NA
TransferQA	35.8		21.3	60.8
UNO-DST-JT	36.6 (+0.8)		36.9 (+15)	75.9 (+15)
UNO-DST-ST	40.8 (+5.0)		47.4 (+26)	81.8 (+21)

Table 1: Average zero-shot JGA and AGA results on MultiWOZ and SGD. JT/ST stands for joint/self-training and **red** figures calculate the performance increase of UNO-DST (60M) over TransferQA (770M).

MultiWOZ Domains	T5DST*	SD-T5	UNO-DST†	UNO-DST‡	JT	ST
Attraction	30.45	33.92	33.50	36.05	32.86	33.09
Hotel	19.38	19.85	21.04	23.00	22.91	25.66
Restaurant	20.42	20.75	22.36	24.03	29.47	30.99
Taxi	66.32	66.25	65.23	65.03	66.00	65.48
Train	25.60	36.96	38.72	47.95	31.68	48.90
Average	32.44	35.55	36.17	39.21	36.58	40.82

Table 2: Zero-shot JGA results for domains in MultiWOZ. **Bold** indicates the best results, * shows results of our replicated T5DST model, while † and ‡ give “t5-small” and “t5-QA” as our model checkpoints. JT/ST stands for joint/self-training.

tributed in both training and testing distributions over 7, 13 domains, representing 7K, 16K training examples in English, respectively. We use standard means for data pre-processing (Budzianowski et al., 2018) for data pre-processing and follow the MultiWOZ leave-one-out settings for zero-shot training and testing in both datasets (Wu et al., 2019; Rastogi et al., 2020).

Evaluation metrics. The primary metric for DST evaluation is Joint Goal Accuracy (JGA), which compares the set of generated predicted values with the set of ground truth after each turn of conversation and returns the fraction of correct matches (Henderson et al., 2014). Average Goal Accuracy (AGA) calculates the JGA only for active slot types, which is used in the SGD baseline (Rastogi et al., 2020).

Baselines and experiment setup. We use T5 (Raffel et al., 2020) as our baseline model. For a fair analysis, we also compare our results with previous DST benchmarks: TRADE (Wu et al., 2019) and the SGD baseline (Rastogi et al., 2020), and current SOTA models: T5DST (Lin et al., 2021b), TransferQA (Lin et al., 2021a) and SD-T5 (Wang et al., 2022). We adopt the cross-domain settings

SGD Domains	TransferQA		UNO-DST-JT		UNO-DST-ST	
	JGA	AGA	JGA	AGA	JGA	AGA
Flights	3.6	42.9	26.4	75.1	25.3	72.7
RideSharing	31.2	61.7	33.3	64.3	73.5	89.8
Homes	31.7	80.6	16.8	77.6	17.9	76.3
Events	15.6	56.8	11.5	58.0	23.1	71.6
Movies	24.0	56.2	35.5	86.7	52.6	86.7
Services	37.2	75.6	75.1	92.1	77.2	92.4
Travel	14.0	24.2	55.2	76.7	56.4	77.8
Weather	40.3	59.4	93.8	98.0	94.3	98.5
Hotels	13.5	60.1	44.8	85.6	75.9	94.6
RentalCars	10.8	73.8	7.5	72.9	5.4	79.4
Restaurants	16.3	68.9	31.8	74.7	35.9	78.5
Media	30.2	67.5	37.0	69.7	60.0	89.2
Music	8.9	62.4	11.6	54.9	19.1	55.5
Average	21.3	60.8	36.9	75.9	47.4	81.8

Table 3: Zero-shot JGA and AGA results for domains in SGD. **Bold** shows the best results and JT/ST stands for joint/self-training.

for both datasets, experimenting on two checkpoints, “t5-small”³ and “t5-QA”⁴. We detail our experimental setup and parameters in § A.1.

5 Results

Table 1 analyzes the zero-shot results of our UNO-DST and the baseline models, including the state-of-the-art (SOTA) TransferQA model, across both datasets. Our model surpasses all baselines for both joint and self-training phases, inclusive of TransferQA using “t5-large” (a magnitude larger in model size). Detailed outcomes for each domain and specific training period (either joint or self-training) across each dataset are presented in Tables 2 and 3.

Joint training results. For the joint training period in the MultiWOZ dataset (Table 2), we are using the same model and prompt as T5DST. UNO-DST shows an increase of more than 4% for JGA across all checkpoints. For SGD (Table 3), our joint training period increases the previous baseline by even larger margins of 15.6% in JGA (15.1% in AGA). The joint training period is critical as it prepares a model for self-training. Table 2 shows that using different model checkpoints is also critical even when using the same model architecture and parameter size. The model performs best when we follow the prompt format in each pre-training checkpoint.

Self-training results. For the MultiWOZ dataset (Table 2), self-training further improves average

³<https://huggingface.co/t5-small>

⁴<https://github.com/facebookresearch/Zero-Shot-DST>

Round	Att.	Hotel	Res.	Taxi	Train	Avg	Gain \uparrow
0	32.86	22.91	29.47	66.00	31.68	36.58	
1	33.09	25.66	30.99	65.48	48.90	40.82	(+4.24)
2	35.53	27.22	31.44	64.71	54.60	42.70	(+1.88)
3	36.62	27.09	31.14	65.48	53.31	42.73	(+0.03)

Table 4: JGA for multiple rounds of self-training on MultiWOZ. Absolute gains indicated in red.

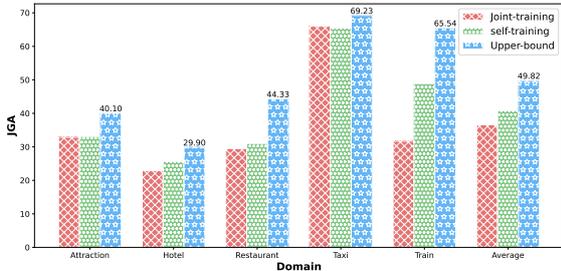


Figure 3: Gains by joint and self-training stages of UNO-DST on the “t5-QA” checkpoint. We show the results of upper-bound (oracle) in each domain for relative comparison.

JGA by 3.09% after joint training. Compared with the baseline, the best performance increases by 8.38% in JGA. In SGD (Table 3), self-training improves the average JGA and AGA in 12 out of 13 domains by an average of 10.5% and 5.9% compared with the joint-training alone, and over 26% and 21% compared to the baseline.

The success of self-training proves the possibility of using pseudo labels generated from zero-shot DST models to bootstrap performance. However, carefully selecting good samples to fine-tune the model is challenging because not all the domains benefit from the self-training process. For example, the result for the “Taxi” domain in MultiWOZ and “Flights” domain in SGD decreased after self-training. We examine the rationale behind the gains obtained through self-training, which is associated with the upper bound results in each domain and further discussed in § 6.

As shown in Table 4, as we lengthen self-training from a single round to multiple rounds, our framework’s performance continues to improve. However, the performance gap between results from different rounds shows diminishing returns, signalling a plateau. The best result with UNO-DST comes when adding more variation is insignificant and so we stop the training when the margin is below 0.1 in JGA. Future work is required to systematically study this strategy over multi-round self-training.

All Generated Slot Types in Train Domain

$people^5, day^5, destination^5, departure^5, leave^5, arrive^5, price^5, type^5, time^5, area^5, name^5$

Valid New Slot Types

$price^{15}, day^1, parking^3, name^5$

Dialogue Example [PMUL1359] ($price^5, name^5$)

System: “Okay, tr6572 departs at 05:29.”

User: “What is the price?”

Dialogue Example [PMUL3027] ($parking^3$)

System: “I have 2 Turkish restaurants in the centre?”

User: “Do they offer free parking?”

Dialogue Example [PMUL1118] (day^1)

User: “I am in Cambridge for the week and want to know what museums you guys have there.”

Table 5: Newly-generated slot types with examples. The superscript on each slot type indicates the domain information from: (1:attraction, 2:hotel, 3:restaurant, 4:taxi, 5:train)

6 Discussion

Lower and upper bound. We calculate the lower and upper bound for zero-shot DST (cf § 3.5) for MultiWOZ. The lower bound is an average of 27.96% for JGA across all domains, a threshold below which the results are unreasonably low (Budzianowski et al., 2018). Of all the baselines in this experiment, only results from TRADE are below the lower bound (Table 1) and based on that, we claim that TRADE (Wu et al., 2019) model is unsuitable for zero-shot DST.

For the upper bound, oracular calculation, we select only the 100% correct samples from the zero-shot predictions and use them for fine-tuning. In Figure 3, we visualise the gains of joint training and self-training alongside our upper bound. While efficacy differs from domain to domain, an important observation is that when the margin between the upper bound (blue columns) and joint training (red columns) is large, the model has a larger gain from self-training, as in “Train” domain. In contrast, for “Taxi” domain, the influence of self-training is weak (cf § 5). Utilizing upper bounds calculations enables us to swiftly evaluate whether a domain or model is apt for the self-training period. Said another way, a larger margin between joint training and the upper bound yields a larger potential improvement that the model can achieve with fine-tuning.

New slot type generation. All the existing zero-

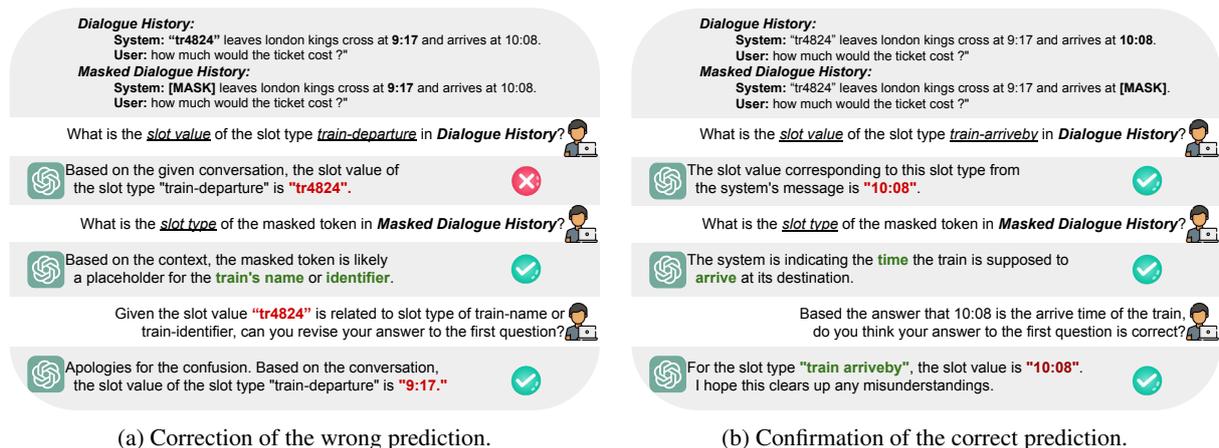


Figure 4: Case studies for conversations with ChatGPT for zero-shot DST.

shot DST methods require given slot types in generating the slot values for both source and target domains and our model also follows the same experiment settings (cf § 3.1). However, we show that our model can also self-generate reasonable slot types either in or beyond the 30 given slot types with our designed auxiliary task. In order to self-generate new slot types, we do a case study on the MultiWoZ “Train” domain and perform random word masking for all the dialogue history, inputting those randomly masked dialogue histories to the auxiliary task for slot type predictions, as shown in Figure 5b. In Table 5, we show some valid new slot types generated by our auxiliary tasks with dialogue history. For example, “asking for the ticket price” in “Train” domain and “asking for parking information” in “Restaurant” domain are reasonable new slot types, which can also be included in the given slot types list.

Besides generating new slot types, we propose a future zero-shot DST without any given slot types, where the self-training is conducted based on self-generated slot types in the target domain. Due to the limited scope of this paper, we discuss the details in Appendix § A.3. Even though our current auxiliary task can predict all 6 given slot types in “Train” domain, additional steps for selection and merging are required to maintain the quality and efficiency of self-generated slot types (Hudeček et al., 2021).

7 UNO-DST with ChatGPT

While earlier sections analyse the effectiveness of our methods on PLMs like “T5”, this section focuses on the potential application of UNO-DST with current generation large language models

(LLMs), such as Brown et al., 2020 and Touvron et al., 2023. Specifically, we examine UNO-DST for zero-shot DST using OpenAI’s ChatGPT⁵ as the backbone LLM, an LLM which has been adopted as a language tool for information extraction with strong capabilities, even without specific training or fine-tuning.

ChatGPT can be interacted with either an open-accessed web interface or through its API. Our study will test the efficacy of our self-training strategy in UNO-DST, on both these versions of ChatGPT, including the web interface and the API⁶. Our objective is to explore the potential application of our methods to LLMs featuring conversational approaches (§7.1) and in-context learning (ICL; Hu et al., 2022 §7.2).

7.1 Conversational approaches

Implementation. We skip the joint training for ChatGPT and use conversations as an inference approach. As shown in Figure 4, we implement main and auxiliary tasks with conversations asking for slot values or types. The selection and fine-tuning steps in the self-training strategy have been converted into the correction or confirmation step, where we provide the slot value and type predictions from the previous two questions to ChatGPT and ask whether it needs to revise the slot value to the main task. We consider the revised response from ChatGPT as the final answer to our main task. We manually examine all the responses generated and give examples of their performance.

Results and discussion. We show two cases of predictions made by ChatGPT using the same

⁵<https://chatgpt.openai.com>

⁶ChatGPT API model: gpt-3.5-turbo-0301

dialogue history in Figure 4a and 4b. In Figure 4a, ChatGPT first made a wrong main task prediction, followed by a correct but not consistent auxiliary task prediction. When we provide all historical information to ChatGPT and ask for a revised main task prediction, it realised that the slot value from the first prediction is not consistent with the original slot type and it self-corrected its wrong prediction. In Figure 4b, we show another case of correct predictions which happens for the majority of the conversations. When ChatGPT is able to make correct predictions for both the main and auxiliary tasks, it confirms the correct predictions for the final question based on cycle consistency.

We illustrate how the cycle consistency between the main and auxiliary tasks aids ChatGPT in rectifying incorrect answers or confirming correct answers provided by the LLM (Zhu et al., 2017). The strategy is applicable to the free accessible web interface of ChatGPT and is easy to implement. However, as conversations are difficult to quantify and evaluate, we only show the results in a qualitative way and encourage future research to explore all the potential implementations in LLMs (Wang et al., 2023).

7.2 In-context Learning

Implementation. Following the settings in § 7.1, we skip the joint training period and apply our strategy to ChatGPT by ICL, providing examples and instructions in the prompt context. Specifically, we compare the ICL results using prompts from two different types of examples: 1) examples from source domain and 2) examples generated and selected by cycle consistency in target domain. To illustrate, the ICL prompt is originally built by DST examples from the source domains and DST test turn from the target domain. After selecting good samples with the strategy discussed in § 7.1, the ICL prompt can be updated with examples in the target domain. We conduct small-scale experiments 3 times with the “train” domain in the MultiWOZ dataset by randomly sampling 100 turns of conversations and inference with ICL prompts using examples generated from the source or target domains, evaluated by JGA.

Results and discussion. The resulting average JGA for the original source domain ICL prompt is 34.92% while the JGA for the selected target domain ICL prompt is 54.18%. Our self-training strategy works very well, serving the LLM to select

valuable in-domain examples for the ICL prompt, improving the zero-shot DST performance by a large 19.25% margin. By manually examining the generated dialogue states, the ICL prompt modified with our strategy performs better, especially for conversations with longer dialogue turns and more slot types. Our self-training strategy demonstrates its capability in generating and selecting dialogue state samples in LLMs which can further improve the performance of zero-shot DST using an ICL prompt. However, due to the scope of this paper, we only test the application of our UNO-DST on ChatGPT with a small corpus of dialogues for case studies (§ 8).

In summary, this section extends the applicability of the UNO-DST strategy to LLMs like ChatGPT, assessing its efficacy in both conversational approaches and ICL. Conversational methods offer a straightforward mechanism for rectifying in-discussion errors, whereas ICL, leveraging APIs or LLM inferences, facilitates handling larger data corpora. Besides testing the cycle consistency strategy in well-suited DST tasks, additional work is required to extend it to other LLMs or more general NLP problems.

8 Conclusion

We propose a novel approach to convert the zero-shot DST into a few-shot setting by generating and selecting quality dialogue states from unlabeled data in the target domain through joint and self-training periods. We introduce and demonstrate how our proposed auxiliary task, which generates slot types as the inverse prompt for the main task which generates slot values, serves the whole model for 1) better accuracy of the main task in joint training 2) quality data selection in the self-training period 3) new slot types generation beyond the given slot type list and 4) upgrading to LLMs without pre-training.

Our proposed strategies of UNO-DST are task-independent, which can be extended to other prompt formats and generalised to LLMs. We look forward to future works that engage additional auxiliary tasks which target new datasets and apply zero-shot DST, even where no slot types are given.

Limitations

This work has two main limitations: 1) due to the limitation of computational resources, we only conduct experiments on small PLMs, which is “t5-

small” and simple NLP tasks, which is “QA”. Our future works will include more NLP tasks with different PLMs to systematically test the performance of our proposed models. 2) Our reported self-training results are only for a single round of self-training because we could not find a way to continuously increase the performance of self-training. Our future plan seeks to improve and examine the best criteria for self-training using an early-stopping strategy. 3) The experimental settings for ChatGPT can be improved in three aspects: a) a larger data corpus can be applied with better instruction prompts in order to limit ChatGPT in generating more accurate values, b) an open-source LLM (Touvron et al., 2023) can be applied to better evaluate and replicate the results of the experiment, and c) a full self-training strategy that includes generation, selection, and fine-tuning can be tested in order to demonstrate the best performance with LLMs.

Ethical Concerns

Our self- and joint-training tunes models to attenuate and amplify signals from the original dataset. While this strategy does work well in our experiments, if the dataset’s signal is weak to start with, our methods may incorrectly amplify errors or biases. Application of our techniques in practical settings should be evaluated before deployment.

This work experimented with publicly available datasets which require no additional annotation from human annotators.

References

Vevake Balaraman, Seyedmostafa Sheikhalishahi, and Bernardo Magnini. 2021. [Recent neural methods on dialogue state tracking for task-oriented dialogue systems: A survey](#). In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 239–251. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*,

volume 33, pages 1877–1901. Curran Associates, Inc.

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, 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*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.

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

Raghav Gupta, Harrison Lee, Jeffrey Zhao, Yuan Cao, Abhinav Rastogi, and Yonghui Wu. 2022. [Show, don’t tell: Demonstrations outperform descriptions for schema-guided task-oriented dialogue](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4541–4549, Seattle, United States. Association for Computational Linguistics.

Michael Heck, Carel van Niekerk, Nurul Lubis, Christian Geishauser, Hsien-Chin Lin, Marco Moresi, and Milica Gašić. 2020. [Trippy: A triple copy strategy for value independent neural dialog state tracking](#). *arXiv preprint arXiv:2005.02877*.

Matthew Henderson, Blaise Thomson, and Jason D. Williams. 2014. [The second dialog state tracking challenge](#). In *Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, pages 263–272, Philadelphia, PA, U.S.A. Association for Computational Linguistics.

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

Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu, Noah A. Smith, and Mari Ostendorf. 2022. [In-context learning for few-shot dialogue state tracking](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2627–2643, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Vojtěch Hudeček, Ondřej Dušek, and Zhou Yu. 2021. [Discovering dialogue slots with weak supervision](#). In *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)*, pages 2430–2442. Association for Computational Linguistics.

702	Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2021.	Hongjie Shi, Takashi Ushio, Mitsuru Endo, Katsuyoshi	760
703	Dialogue state tracking with a language model using	Yamagami, and Noriaki Horii. 2017. Convolutional	761
704	schema-driven prompting . In <i>Proceedings of the</i>	neural networks for multi-topic dialog state track-	762
705	<i>2021 Conference on Empirical Methods in Natural</i>	ing. <i>Dialogues with Social Robots: Enablements,</i>	763
706	<i>Language Processing</i> , pages 4937–4949, Online and	<i>Analyses, and Evaluation</i> , pages 451–463.	764
707	Punta Cana, Dominican Republic. Association for		
708	Computational Linguistics.		
709	Shuyang Li, Jin Cao, Mukund Sridhar, Henghui Zhu,	Jamin Shin, Hangyeol Yu, Hyeongdon Moon, Andrea	765
710	Shang-Wen Li, Wael Hamza, and Julian McAuley.	Madotto, and Juneyoung Park. 2022. Dialogue sum-	766
711	2021. Zero-shot generalization in dialog state track-	maries as dialogue states (DS2), template-guided	767
712	ing through generative question answering . In <i>Pro-</i>	summarization for few-shot dialogue state tracking .	768
713	<i>ceedings of the 16th Conference of the European</i>	In <i>Findings of the Association for Computational</i>	769
714	<i>Chapter of the Association for Computational Lin-</i>	<i>Linguistics: ACL 2022</i> , pages 3824–3846, Dublin,	770
715	<i>guistics: Main Volume</i> , pages 1063–1074, Online.	Ireland. Association for Computational Linguistics.	771
716	Association for Computational Linguistics.		
717	Zhaojiang Lin, Bing Liu, Andrea Madotto, Seungwhan	Yixuan Su, Lei Shu, Elman Mansimov, Arshit Gupta,	772
718	Moon, Zhenpeng Zhou, Paul Crook, Zhiguang Wang,	Deng Cai, Yi-An Lai, and Yi Zhang. 2022. Multi-task	773
719	Zhou Yu, Eunjoon Cho, Rajen Subba, and Pascale	pre-training for plug-and-play task-oriented dialogue	774
720	Fung. 2021a. Zero-shot dialogue state tracking via	system . In <i>Proceedings of the 60th Annual Meet-</i>	775
721	cross-task transfer . In <i>Proceedings of the 2021 Con-</i>	<i>ing of the Association for Computational Linguistics</i>	776
722	<i>ference on Empirical Methods in Natural Language</i>	<i>(Volume 1: Long Papers)</i> , pages 4661–4676, Dublin,	777
723	<i>Processing</i> , pages 7890–7900, Online and Punta	Ireland. Association for Computational Linguistics.	778
724	Cana, Dominican Republic. Association for Com-		
725	putational Linguistics.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	779
726	Zhaojiang Lin, Bing Liu, Seungwhan Moon, Paul	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	780
727	Crook, Zhenpeng Zhou, Zhiguang Wang, Zhou Yu,	Baptiste Rozière, Naman Goyal, Eric Hambro,	781
728	Andrea Madotto, Eunjoon Cho, and Rajen Subba.	Faisal Azhar, et al. 2023. Llama: Open and effi-	782
729	2021b. Leveraging slot descriptions for zero-shot	cient foundation language models. <i>arXiv preprint</i>	783
730	cross-domain dialogue StateTracking . In <i>Proce-</i>	<i>arXiv:2302.13971</i> .	784
731	<i>eedings of the 2021 Conference of the North Ameri-</i>		
732	<i>can Chapter of the Association for Computational</i>	Yao-Hung Hubert Tsai, Yue Wu, Ruslan Salakhutdinov,	785
733	<i>Linguistics: Human Language Technologies</i> , pages	and Louis-Philippe Morency. 2021. Self-supervised	786
734	5640–5648, Online. Association for Computational	learning from a multi-view perspective .	787
735	Linguistics.		
736	Nikola Mrkšić, Diarmuid Ó Séaghdha, Tsung-Hsien	Qingyue Wang, Yanan Cao, Piji Li, Yanhe Fu, Zheng	788
737	Wen, Blaise Thomson, and Steve Young. 2017. Neu-	Lin, and Li Guo. 2022. Slot dependency model-	789
738	ral belief tracker: Data-driven dialogue state tracking .	ing for zero-shot cross-domain dialogue state track-	790
739	In <i>Proceedings of the 55th Annual Meeting of the</i>	ing . In <i>Proceedings of the 29th International Con-</i>	791
740	<i>Association for Computational Linguistics (Volume 1:</i>	<i>ference on Computational Linguistics</i> , pages 510–520,	792
741	<i>Long Papers)</i> , pages 1777–1788, Vancouver, Canada.	Gyeongju, Republic of Korea. International Commit-	793
742	Association for Computational Linguistics.	tee on Computational Linguistics.	794
743	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	795
744	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,	796
745	Wei Li, and Peter J. Liu. 2020. Exploring the limits	and Denny Zhou. 2023. Self-consistency improves	797
746	of transfer learning with a unified text-to-text trans-	chain of thought reasoning in language models . In	798
747	former. <i>J. Mach. Learn. Res.</i> , 21(1).	<i>The Eleventh International Conference on Learning</i>	799
748	Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara,	<i>Representations</i> .	800
749	Raghav Gupta, and Pranav Khaitan. 2020. Towards	Chien-Sheng Wu, Steven C.H. Hoi, Richard Socher,	801
750	scalable multi-domain conversational agents: The	and Caiming Xiong. 2020a. TOD-BERT: Pre-trained	802
751	schema-guided dialogue dataset. In <i>Proceedings of</i>	natural language understanding for task-oriented di-	803
752	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	alogue . In <i>Proceedings of the 2020 Conference on</i>	804
753	ume 34, pages 8689–8696.	<i>Empirical Methods in Natural Language Processing</i>	805
754	Liliang Ren, Kaige Xie, Lu Chen, and Kai Yu. 2018. To-	<i>(EMNLP)</i> , pages 917–929, Online. Association for	806
755	wards universal dialogue state tracking . In <i>Proce-</i>	Computational Linguistics.	807
756	<i>eedings of the 2018 Conference on Empirical Methods</i>	Chien-Sheng Wu, Steven C.H. Hoi, and Caiming Xiong.	808
757	<i>in Natural Language Processing</i> , pages 2780–2786,	2020b. Improving limited labeled dialogue state	809
758	Brussels, Belgium. Association for Computational	tracking with self-supervision . In <i>Findings of the</i>	810
759	Linguistics.	<i>Association for Computational Linguistics: EMNLP</i>	811
		2020, pages 4462–4472, Online. Association for	812
		Computational Linguistics.	813
		Chien-Sheng Wu, Andrea Madotto, Ehsan Hosseini-Asl,	814
		Caiming Xiong, Richard Socher, and Pascale Fung.	815

2019. Transferable multi-domain state generator for task-oriented dialogue systems. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 808–819, Florence, Italy. Association for Computational Linguistics.

Yuting Yang, Wenqiang Lei, Pei Huang, Juan Cao, Jintao Li, and Tat-Seng Chua. 2023. A dual prompt learning framework for few-shot dialogue state tracking. In *Proceedings of the ACM Web Conference 2023, WWW '23*, page 1468–1477, New York, NY, USA. Association for Computing Machinery.

Yu Zhang and Qiang Yang. 2021. A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering*, 34(12):5586–5609.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Computer Vision (ICCV), 2017 IEEE International Conference on*.

A Appendix

A.1 Experimental Setup

We select QA as the main task in our framework for its popularity and test it through different checkpoints holding parameters fixed. We adopt the open-source T5-SMALL (Raffel et al., 2020) with 60M parameters as our baseline and train using AdamW with a learning rate of 0.0001, batch size 8 for 1 epoch (zero-shot setting) and 3 epochs (fine-tuning setting). We initialize the model using “t5-small”⁷ and “t5-QA”⁸ checkpoints. For “t5-small”, the same question prompt is used as in T5DST (Lin et al., 2021b). While for “t5-QA” which is pre-trained following the steps in *TransferQA*, a similar prompt format is chosen (Lin et al., 2021a). We train on a single GeForce RTX3000 GPU.

During the training period, we adopt the cross-domain setting for MultiWOZ 2.1 all domains (Wu et al., 2020a) and SGD seen domains (Rastogi et al., 2020). For example, under this setting, we use four of the five domains in MultiWOZ as the source domains for the training of the model. The remaining domain is used as the target domain for testing performance (Wu et al., 2019). The unlabelled training data in the target domain will be used for self-training and testing data in target domain is only used for final testing.

Algorithm 1 Self-training with auxiliary task

Require: Training dataset K_u in unknown domains, slot type $s \in S$, dialogue history C_t and PLM $PLM(\cdot)$ after joint training
Ensure: The fine-tuned model $PLM(\cdot)$

```

1: repeat
2:   Step 1
3:   for batch  $k_u \in K_u$  do
4:     for slot-type, context  $s, C_t \in k_u$  do
5:        $p^{main} = P(s, C_t) \rightarrow$ prompting
6:        $v' \leftarrow PLM(p^{main}) \rightarrow$ main task
7:        $C_t^m = M(C_t, v') \rightarrow$ value masking
8:        $p^{aux} = IP(C_t^m) \rightarrow$ inverse prompting
9:        $s' \leftarrow PLM(p^{aux}) \rightarrow$ auxiliary task
10:    end for
11:  end for
12:  Step 2
13:   $G = [ ]$ 
14:  for batch  $k_u \in K_u$  do
15:    if  $set(s) = set(s')$  then
16:       $G.append(k_u)$ 
17:    end if
18:  end for
19:  Step 3 fine-tuning  $PLM(\cdot)$  with  $G$ 

```

A.2 Self-training Algorithm

In this section, we explain the detail for our self-training strategy in Algorithm 1. The self-training is conducted using PLM after the joint training and based on the unlabelled training data in the target domain. There are a total of 3 steps for self-training which are 1) generation of slot values and types from the main and auxiliary task, 2) selection of good dialogue state samples using cycle consistency between two tasks and 3) Fine-tuning the PLM with selected samples.

A.3 Slot type discussion

A.3.1 Future zero-shot DST

As discussed in the previous section, the auxiliary task can facilitate the generation of new slot types beyond the pre-defined 30 slots in MultiWOZ. Besides adding more slot types to the given slot type list, we believe that it is possible for our proposed model to conduct DST tasks in an unknown target domain without any given slot types and we describe the proposal of future zero-shot DST in

⁷<https://huggingface.co/t5-small>

⁸<https://github.com/facebookresearch/Zero-Shot-DST/tree/main/TransferQA>

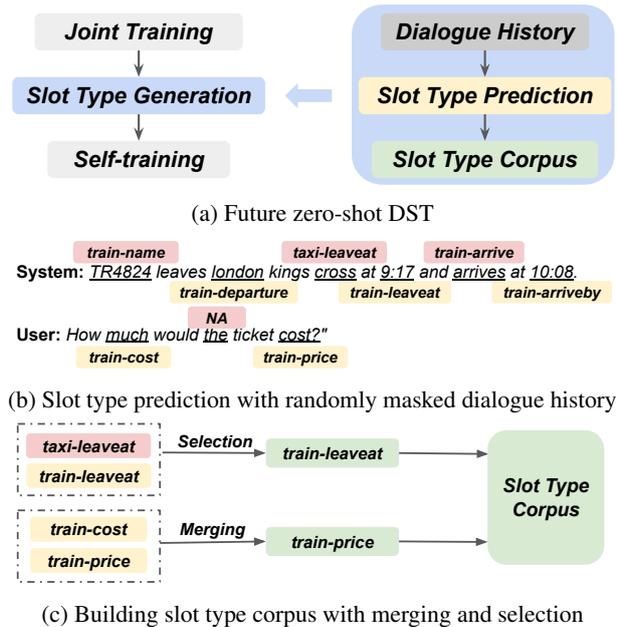


Figure 5: Zero-shot DST without pre-defined slot types

Figure 5

In order to eliminate the use of pre-defined slot types, we add a slot type generation period between joint and self-training, which identifies and select domain-relevant slot type corpus. Similar to the process proposed by Hudeček et al. (2021), we can first use our auxiliary task to generate potential slot types based on random masked dialogue history, as shown in Figure 5b. The generated text may contain domain-irrelevant or similar slot types and we propose a weak selection and merging of task-relevant and similar slot types for slot type corpus (Hudeček et al., 2021). Secondly, those generated slot-type corpus can be used for self-training in the unknown target domain, as discussed in § 3.4.

During our testing, our auxiliary task can generate predictions including all 6 given slot types in the “train” domain, as well as valid slot types in other domains, which demonstrates the potential of future zero-shot methods without pre-defined slot types. We look forward to future works for zero-shot DST without any labelled data in slot types and values in unknown target domains.

A.3.2 Unseen slot type prediction

For each domain, there are seen slot types which appear in other domains (e.g., “hotel-name” and “restaurant-name”) or unseen slot types which are unique in the specific domain (e.g., “hotel-stars”). The ratio of the occurrences of these slot types represents the difficulty of zero-shot DST for each

domain (Wang et al., 2022). As shown in table 6, the original setting for the MultiWoz dataset has 30 given slot types. However, not all of them appear in every domain. For some certain domains, like the “hotel” domain, there are 4 unique slot types which do not appear in other domains, including “stars”, “internet”, “stay” and “parking”. In Figure 6, we show the slot accuracy for the hotel domain. It shows that generally, the unseen slot types will perform worse than the seen slot types (Wang et al., 2022). Prediction for those slot types in zero-shot cross-domain settings can be challenging as there is no further information from the other source domains. In addition, half of the unseen slot types in the “hotel” domain are related to “yes/no” slot values, whereas in our joint training settings in § 3.3, we skip the masking of those “yes/no” values from the context and the model is less trained compared with other slot types. We hope that future works should improve on “yes” or “no” value prediction.

All Given Slot Types in MultiWOZ 2.1

*area*¹²³, *arriveby*⁴⁵, *day*²³⁵, *departure*⁴⁵, *destination*⁴⁵, *food*³, *internet*², *leave*⁴⁵, *name*¹²³, *people*²³⁵, *parking*², *price*²³, *stars*², *stay*², *time*³, *type*¹²

Seen Slot Types in Hotel Domain

*area*¹²³, *day*²³⁵, *name*¹²³, *people*²³⁵, *price*²³, *type*¹²

Unseen Slot Types in Hotel Domain

*internet*², *parking*², *stars*², *stay*²

Table 6: Seen and unseen slot types in hotel domain. The superscript on each slot type indicates the domain information from: (1:attraction, 2:hotel, 3:restaurant, 4:taxi, 5:train)

A.3.3 Auxiliary Task Design

In our model, in order to better facilitate the main task with cycle consistency, we design the auxiliary task as an inverse prompt of the main task and in order to better leverage the capability PLMs, we convert the auxiliary prompt as mask language prediction. We conclude our design of the auxiliary task into 3 criteria: 1) The auxiliary task and the main task should have a similar format so that they can share the prompt function P . 2) The auxiliary task needs to facilitate the main task in one or many different ways. In UNO-DST, the auxiliary task serves the main task in two different ways, which are generation and selection. 3) The auxiliary task should be easier than the main task so

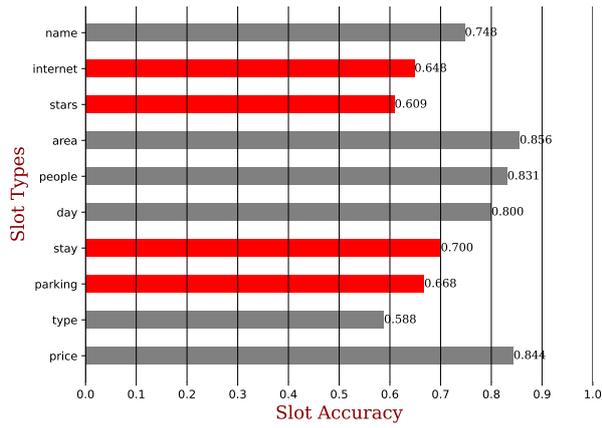


Figure 6: Slot Accuracy for seen and unseen slot types in hotel domain. (Grey: Seen slot types, Red: Unseen slot types)

948 that it can better serve the selection in self-training.
 949 Our work shows the initiative of adding auxiliary
 950 tasks and we hope that future work can propose
 951 different tasks, such as targeting “none” or “yes/no”
 952 value predictions. Besides cycle consistency, the
 953 self-consistency between multiple auxiliary tasks
 954 should also help with the main task in generation
 955 and selection (Wang et al., 2023).