

Schema Encoding for Transferable Dialogue State Tracking

Anonymous ACL submission

Abstract

Dialogue state tracking (DST) is an essential sub-task for task-oriented dialogue systems. Recent work has focused on deep neural models for DST. However, the neural models require a large dataset for training. Furthermore, applying them to another domain needs a new dataset because the neural models are trained to imitate the given dataset. In this paper, we propose Schema Encoding for Transferable Dialogue State Tracking (SET-DST), which is a neural DST method for effective transfer to new domains. Transferable DST could assist developments of dialogue systems even with few dataset on target domains. We use a schema encoder not just to imitate the dataset but to comprehend the schema of the dataset. We aim to transfer the model to new domains by encoding new schemas and using them for DST. As a result, SET-DST improved the accuracy by 1.46 points on MultiWOZ 2.1.

1 Introduction

The objective of task-oriented dialogue systems is to help users achieve their goals by conversations. Dialogue state tracking (DST) is the essential sub-task for the systems to perform the purpose. Users may deliver the details of their goals to the systems during the conversations, e.g., what kind of food they want the restaurant to serve and at what price level they want to book the hotel. Thus, the systems should exactly catch the details from utterances. They should also communicate with other systems by using APIs to achieve users' goals, e.g., to search restaurants and to reserve hotels. The goal of DST is to classify the users' intents and to fill the details into predefined templates that are used to call APIs.

Recent work has used deep neural networks for DST with supervised learning. They have improved the accuracy of DST; however, they require a large dataset for training. Furthermore, they need

a new dataset to be trained on another domain. Unfortunately, the large dataset for training a DST model is not easy to be developed in real world. The motivation of supervised learning is to make deep neural networks imitate humans. But, they actually imitate the given datasets rather than humans. Someones who have performed hotel reservation work can easily perform restaurant reservation work if some guidelines are provided, but neural models may have to be trained on a new dataset of the restaurant domain. The difference between humans and neural models is that humans can learn how to read guidelines and to apply the guidelines to their work. This is why transfer learning is important to train neural models on new domains.

In this paper, we propose Schema Encoding for Transferable Dialogue State Tracking (SET-DST), which is a neural DST method with transfer learning by using dataset schemas as guidelines for DST. The motivation of this study is that humans can learn not only how to do their work, but also how to apply the guidelines to the work. We aim to make a neural model learn how to apply the schema guidelines to DST beyond how to fill predefined slots by simply imitating the dataset. The schema includes metadata of the dataset, e.g., which domains the dataset covers and which slots have to be filled to achieve goals. SET-DST has a schema encoder to represent the dataset schema, and it uses the schema representation to understand utterances and to fill slots. Recently, transfer learning has been becoming important because development of new datasets is costly. Transfer learning makes it possible to pre-train neural models on large-scale datasets to effectively fine-tune the models on small-scale downstream tasks.

We evaluated SET-DST on MultiWOZ 2.1 (Eric et al., 2020), which is a standard benchmark dataset for DST, as a downstream task. SET-DST achieved state-of-the-art accuracy on the DST task. We further confirmed that SET-DST worked well when

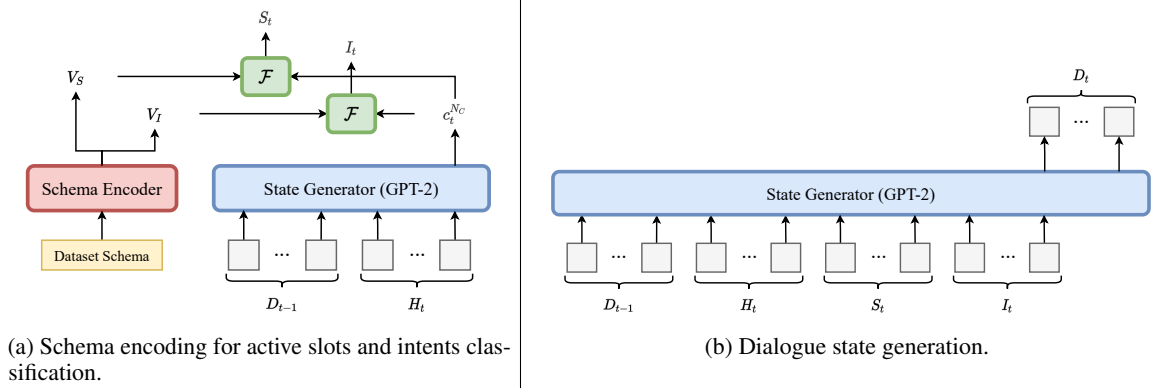


Figure 1: Overview of SET-DST. The schema encoder takes the dataset schema and generates slot vectors and intent vectors. The state generator takes the previous dialogue state D_{t-1} and the dialogue history H_t to calculate active scores of slots and intents. \mathcal{F} is an score function to calculate whether the slots or intents are activated on turn t . Then, the state generator additionally takes the activated slots and intents to generate the current dialogue state D_t . S_t indicates the activated slots and I_t indicates the activated intents.

082 the size of the downstream dataset was small.
 083 This result demonstrates that transfer learning with
 084 schema encoding improves the performance of neural
 085 DST models and the efficiency of few-shot
 086 learning on DST.

087 2 Related Work

088 Traditional DST models extract semantics by using
 089 natural language understanding (NLU) modules to
 090 generate dialogue states (Williams, 2014; Wang
 091 and Lemon, 2013). The limitation of these models
 092 is that they rely on features extracted by humans.

093 Recent work has focused on building end-to-end
 094 DST models without hand-crafted features. Zhong
 095 et al. (2018) use global modules to share parameters
 096 between different slots. Nouri and Hosseini-Asl
 097 (2018) improve the latency by removing inefficient
 098 recurrent layers. Transferable DST models that
 099 can be adapted to new domains by removing the
 100 dependency on the domain ontology are proposed
 101 (Ren et al., 2018; Wu et al., 2019). Zhou and Small
 102 (2019) attempt to solve DST as a question answering
 103 task using knowledge graph.

104 More recently, large-scale pre-trained language
 105 models such as BERT (Devlin et al., 2019) and
 106 GPT-2 (Radford et al., 2019) are used for DST. The
 107 pre-trained BERT acts as an NLU module to under-
 108 stand utterances (Lee et al., 2019; Zhang et al.,
 109 2020a; Kim et al., 2020; Heck et al., 2020). GPT-2
 110 makes it possible to solve DST as a conditional
 111 language modeling task (Hosseini-Asl et al., 2020;
 112 Peng et al., 2021).

113 Rastogi et al. (2020) propose the baseline
 114 method that defines the schema of dataset and uses

115 it for training and inference. A drawback of them
 116 is that the calculation cost is high because they use
 117 the domain ontology and access all values to es-
 118 timate the dialogue state. DST models that uses
 119 schema graphs to encode the relation between slots
 120 and values are proposed (Chen et al., 2020; Zhu
 121 et al., 2020). However, they focus on encoding
 122 the relation between slots and values of the given
 123 domains not on adaptation to new domains.

124 In this paper, we focus on making the model
 125 learn how to understand the schema and how to
 126 apply it to estimate the dialogue state, not just on
 127 encoding the in-domain relation.

128 3 Schema Encoding for Transferable 129 Dialogue State Tracking

130 In this section, we describe the architecture of SET-
 131 DST and how to optimize it. Figure 1 shows the
 132 overview of our method. The model consists of
 133 the schema encoder and the state generator. SET-
 134 DST generates the dialogue state in two steps: (a)
 135 schema encoding and classification and (b) dia-
 136 logue state generation. In this paper, we define
 137 some terms as follows.

138 **Schema** What domains, services, slots, and in-
 139 tents the dataset covers. A dataset has a schema
 140 that describes the dataset.

141 **Domain** What domains the conversation goes on,
 142 e.g., restaurant, hotel, and attraction. A conversa-
 143 tion can go on multiple domains.

144 **Service** What services the system provides to
 145 users. It is similar to domain, but smaller and

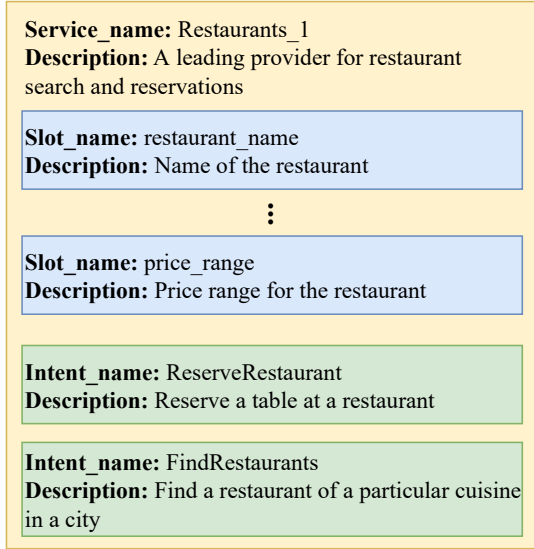


Figure 2: Example of schema for restaurant search and reservation service including slots and intents.

application-level. For example, restaurant domain can have two different services: (1) a service for searching and reserving restaurants and (2) a service focused on searching and comparing restaurants. In real world, a service corresponds to an application.

Action Abstract actions of users to achieve their goals during conversations, e.g., to inform the system their requirements or to request the system for some information.

Slot The details of the user goals, e.g., the type of food and the price range of hotel. Slots are predefined based on the domains or services that the system should cover, and the slots are filled by DST.

Value The values that have actual meaning for corresponding slots, e.g., cheap or expensive about the price range of hotel. The systems should match slot-value pairs from conversations.

Intent Sub-goals to achieve the user goals by conversation. A goal consists of one or more intents, and an intent is achieved over one or more conversation turns. In real world, an intent corresponds to an API. For example, to search restaurants or to book hotels should be performed by APIs of external systems.

3.1 Schema Encoding

We use the pre-trained BERT¹ for the schema encoder. Figure 2 shows an example of the schema

for *Restaurant_1* service that is a service to search and reserve restaurants. Services, slots, and intents consist of name and short description. The name and description of the service in the schema are fed into BERT to generate service vector v_R as

$$\begin{aligned} o_R &= \text{BERT}([\text{CLS}]n_R : d_R[\text{SEP}]) \\ v_R &= W_R \cdot o_R^{[\text{CLS}]} \in \mathbb{R}^h \end{aligned}, \quad (1)$$

where n_R is the service name, d_R is the service description, and h is the hidden size. $o_R^{[\text{CLS}]}$ is the output of [CLS] token, and $W_R \in \mathbb{R}^{h \times h}$ is a fully connected (FC) layer. [CLS] and [SEP] are special tokens that mean the start and end of the sentence, respectively. The service in Figure 2 can be represented as [CLS] Restaurants_1 : A leading provider for restaurant search and reservations [SEP]. The slots and intents in the schema are also fed into BERT to generate slot vectors $V_S = \{v_S^1, \dots, v_S^{N_S}\} \in \mathbb{R}^{N_S \times h}$ and intent vectors $V_I = \{v_I^1, \dots, v_I^{N_I}\} \in \mathbb{R}^{N_I \times h}$, respectively, as follows:

$$\begin{aligned} o_S^j &= \text{BERT}([\text{CLS}]n_S^j : d_S^j[\text{SEP}]) \\ v_S^j &= W_S \cdot o_S^{j, [\text{CLS}]} \in \mathbb{R}^h, \quad j \in [1, N_S] \\ o_I^k &= \text{BERT}([\text{CLS}]n_I^k : d_I^k[\text{SEP}]) \\ v_I^k &= W_I \cdot o_I^{k, [\text{CLS}]} \in \mathbb{R}^h, \quad k \in [1, N_I] \end{aligned}, \quad (2)$$

N_S and N_I mean the number of slots and intents for the service, respectively. n_S^j is the j -th slot name, and d_S^j is the j -th slot description. $o_S^{j, [\text{CLS}]}$ is the output of [CLS] token from the j -th slot, and $W_S \in \mathbb{R}^{h \times h}$ is an FC layer. Similarly, n_I^k is the k -th intent name, and d_I^k is the k -th intent description. $o_I^{k, [\text{CLS}]}$ is the output of [CLS] token from the k -th intent, and $W_I \in \mathbb{R}^{h \times h}$ is an FC layer. The schema encoder takes v_R , V_S , and V_I to update the slot vectors V_S and intent vectors V_I with attention mechanism as follows:

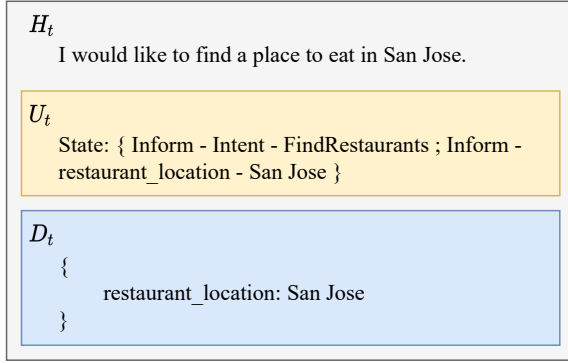
$$\begin{aligned} a_S &= \text{softmax}(V_S \cdot v_R) \in \mathbb{R}^{N_S} \\ v_{R,S} &= (V_S)^T \cdot a_S \in \mathbb{R}^h \end{aligned}, \quad (4)$$

$$v_S^j = W_{RS} \cdot (v_{R,S} \oplus v_S^j) \in \mathbb{R}^h, \quad (5)$$

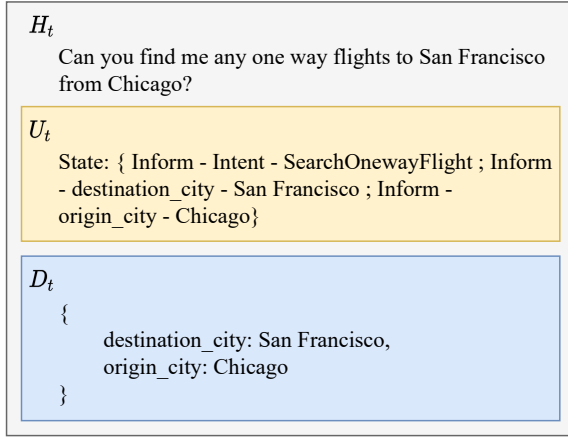
$$\begin{aligned} a_I &= \text{softmax}(V_I \cdot v_R) \in \mathbb{R}^{N_I} \\ v_{R,I} &= (V_I)^T \cdot a_I \in \mathbb{R}^h \end{aligned}, \quad (6)$$

$$v_I^k = W_{RI} \cdot (v_{R,I} \oplus v_I^k) \in \mathbb{R}^h. \quad (7)$$

$W_{RS} \in \mathbb{R}^{h \times 2h}$ and $W_{RI} \in \mathbb{R}^{h \times 2h}$ are FC layers, and \oplus means the concatenation of two vectors.



(a) Example in restaurants domain.



(b) Example in flights domain.

Figure 3: Examples of user state and dialogue state corresponding to user utterance. U_t is a sequence of words, and D_t is a list of slot-value pairs.

3.2 Slot and Intent Classification

SET-DST takes the slot vectors and intent vectors to classify what slots and intents are activated by users. We use the pre-trained GPT-2¹ for the state generator that encodes the dialogue history and generates the dialogue state as a sequence of words. The state generator encodes the dialogue history H_t that is accumulated during the conversation and the previous dialogue state D_{t-1} to calculate the context vector C_t as

$$\begin{aligned}
 C_t &= \{c_t^1, \dots, c_t^{N_C}\} \in \mathbb{R}^{N_C \times h} \\
 &= \text{GPT-2}(D_{t-1} \oplus H_t),
 \end{aligned} \tag{8}$$

where $N_C = |C_t|$. Then, the last output of C_t is used to classify which slots and intents are activated

¹The pre-trained models are available at <https://github.com/huggingface/transformers>.

in the current conversation as follows:

$$P(s_t^j = \text{Active}) = \mathcal{F}(c_t^{N_C}, v_S^j), \tag{9}$$

$$P(i_t^k = \text{Active}) = \mathcal{F}(c_t^{N_C}, v_I^k), \tag{10}$$

where $P(s_t^j = \text{Active})$ means the probability that the j -th slot is activated on turn t , and $P(i_t^k = \text{Active})$ means the probability that the k -th intent is activated on turn t . \mathcal{F} is a projection layer to calculate the probabilities using the context vector, slot vector, and intent vector. We define $P = \mathcal{F}(x, y)$ as a function transforming vectors x and y into a probability scalar as

$$\begin{aligned}
 h_1 &= \tanh(W_1 \cdot x) \\
 h_2 &= \tanh(W_2 \cdot (h_1 \oplus y)), \\
 P &= \sigma(W_3 \cdot h_2)
 \end{aligned} \tag{11}$$

where $W_1 \in \mathbb{R}^{h \times h}$, $W_2 \in \mathbb{R}^{h \times 2h}$, and $W_3 \in \mathbb{R}^{1 \times h}$ are FC layers. Activate slots and intents are classified based on the probabilities $P(s_t^j = \text{Active})$ and $P(i_t^k = \text{Active})$. We define the slots activated on turn t as $S_t = \{s_t^j | P(s_t^j = \text{Active}) \geq \alpha, j \in [1, N_S]\}$ and the intents activated on turn t as $I_t = \{i_t^k | P(i_t^k = \text{Active}) \geq \alpha, k \in [1, N_I]\}$. α is a threshold to classify the slots and intents based on the probabilities.

3.3 Dialogue State Generation

SET-DST has the state generator that generates dialogue state using the dialogue history, schema representation, and previous dialogue state accumulated during the conversation. In this paper, we define the dialogue state as a list of slot-value pairs that mean the details of an user goal. We also define the concept called user state that is a sequence of action-slot-value triples to generalize semantics from various user utterances. The state generator recurrently generates the user state as a sequence of words, and updates the dialogue state by extracting the slot-value pairs from the user state. The user state U_t on turn t is generated based on the previous dialogue state D_{t-1} , dialogue history H_t , active slots S_t , and active intents I_t as follows:

$$\tilde{u}_t^l = \text{GPT-2}(D_{t-1} \oplus H_t \oplus S_t \oplus I_t \oplus U_t^{1:l-1}), \tag{12}$$

$$\begin{aligned}
 U_t &= \{u_t^l | u_t^l = \arg\max(W_{vocab} \cdot \tilde{u}_t^l), \\
 &\quad l \in [1, N_U]\} \in \mathbb{R}^{N_U},
 \end{aligned} \tag{13}$$

where $U_t^{1:l-1} = \{u_t^1, \dots, u_t^{l-1}\}$, $N_U = |U_t|$, and $W_{vocab} \in \mathbb{R}^{N_{vocab} \times h}$ is an FC layer to project the hidden state to vocabulary space with size of N_{vocab} . Figure 3 shows how to generate D_t from U_t . U_t is generated word-by-word over time steps until [EOS], a special word to terminate the generation, is detected. Then, D_t is updated by extracting the slot-value pairs from U_t .

3.4 Optimization

SET-DST is optimized over two steps: (1) slot and intent classification and (2) state generation. We freeze the pre-trained BERT during training to preserve the broad and general knowledge that is learned from large corpus. In classification task, the system is trained by using binary cross-entropy. Equation 9 is used to calculate the slot loss \mathcal{L}_t^S with slot labels $Y_t^S = \{y_{t,1}^S, \dots, y_{t,N_S}^S\}$ as

$$\mathcal{L}_t^S = -\frac{1}{N_S} \sum_{j=1}^{N_S} \beta \cdot y_{t,j}^S \cdot \log P(s_t^j) + (1 - y_{t,j}^S) \log(1 - P(s_t^j)), \quad (14)$$

where $y_{t,j}^S \in \mathbb{R}^1$ is the binary value of j -th slot on turn t , and β is a hyperparameter to consider the ratio of active slots out of total slots. Based on Equation 10, the intent loss \mathcal{L}_t^I is calculated with intent labels $Y_t^I = \{y_{t,1}^I, \dots, y_{t,N_I}^I\}$ as

$$\mathcal{L}_t^I = -\frac{1}{N_I} \sum_{k=1}^{N_I} \beta \cdot y_{t,k}^I \cdot \log P(i_t^k) + (1 - y_{t,k}^I) \log(1 - P(i_t^k)), \quad (15)$$

where $y_{t,k}^I \in \mathbb{R}^1$ is the binary value of k -th intent on turn t . In state generation step, the system is trained as a conditional language model that recurrently generates words over time steps. The state loss \mathcal{L}_t^U is calculated base on Equation 13 with the state label $Y_t^U = \{y_{t,1}^U, \dots, y_{t,N_U}^U\}$ as

$$\mathcal{L}_t^U = -\frac{1}{N_U} \sum_{l=1}^{N_U} (y_{t,l}^U)^T \log P(u_t^l), \quad (16)$$

where $y_{t,l}^U \in \mathbb{R}^{N_{vocab}}$ is the one-hot vector that indicates the l -th word of the gold-standard user state on turn t . The final joint loss is the sum of above losses:

$$\mathcal{L}_t = \mathcal{L}_t^S + \mathcal{L}_t^I + \mathcal{L}_t^U. \quad (17)$$

We use Adam optimizer (Kingma and Ba, 2014) to minimize \mathcal{L}_t .

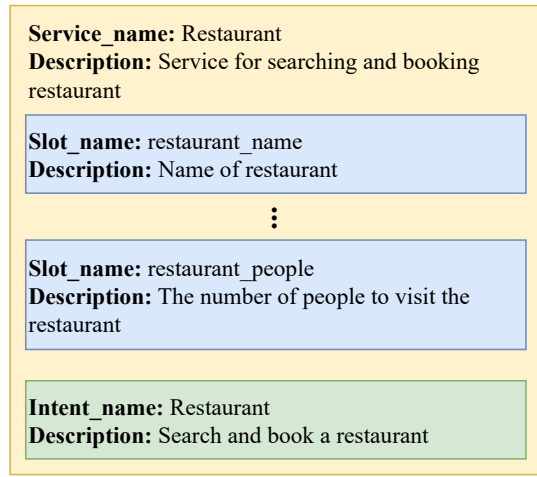


Figure 4: Example of schema that is temporarily created for MultiWOZ dataset.

4 Experiments

In this section, we describe our experiments including the datasets, evaluation metric, and results.

4.1 Experimental Setups

We used two datasets MultiWOZ 2.1² and Schema-Guided Dialogue (SGD)³ to evaluate our system. MultiWOZ consists of conversations between a tourist and a guide, e.g., booking hotels and searching trains. SGD deals with conversations between a virtual assistant and an user ranging over various domains, e.g., events, restaurants, and media. The dataset also provides a schema that includes services, intents, and slots with short descriptions to help understanding the conversations. MultiWOZ has about 10,400 dialogues, and SGD has about 22,800 dialogues.

The datasets propose joint accuracy as the metric to evaluate DST systems. Joint accuracy measures whether a system successfully predicts all slot-value pairs mentioned on the conversations. In every turn, the system updates dialogue state, and the joint accuracy is calculated based on the accumulated dialogue state.

4.2 Experimental Details

The motivation of SET-DST is to make the system interpret the schema and refer it for efficiently tracking the dialogue state. In the experiments, our

²<https://github.com/budzianowski/multiwoz>.

³<https://github.com/google-research-datasets/dstc8-schema-guided-dialogue>.

<p>Service_name</p> <p>Original: Banks</p> <p>Alternatives:</p> <pre>[Bank_service, Bank_application, ...]</pre>
<p>Service_description</p> <p>Original: Manage bank accounts and transfer money</p> <p>Alternatives:</p> <pre>[Service to manage your bank accounts and finances, Application for managing bank accounts, ...]</pre>
<p>Slot_name</p> <p>Original: account_type</p> <p>Alternatives:</p> <pre>[bank_account_type, type_of_bank_account, ...]</pre>
<p>Slot_description</p> <p>Original: The account type of the user</p> <p>Alternatives:</p> <pre>[Bank account type of the user for transaction, Type of user's bank account, ...]</pre>
<p>Intent_name</p> <p>Original: transfer_money</p> <p>Alternatives:</p> <pre>[send_money, money_transference, ...]</pre>
<p>Intent_description</p> <p>Original: Transfer money from one bank account to another user's account</p> <p>Alternatives:</p> <pre>[Transfer money to another user, Send money to another bank account, ...]</pre>

Figure 5: Example of alternatives for schema augmentation.

goal is to verify that SET-DST works well for our purpose by improving the performance of DST and the efficiency on few-shot settings with the schema encoding.

The experiments are divided into two steps: (1) pre-training on SGD and (2) fine-tuning on MultiWOZ. In the pre-training step, SET-DST is optimized to encode the schema for DST. In the fine-tuning step, the capability that encodes given schema is transferred to encode new schema for improvement of the performance and efficiency. We conducted the experiments by adjusting the rate of few-shot data during fine-tuning to focus on the fine-tuning step. The training data for few-shot settings was randomly sampled from the training set of MultiWOZ, and the random seed was fixed

Hidden size	768
Embedding size	768
Vocabulary size	30522
Dropout	0.3
Early stopping count	5
Max epochs	40
Min epochs	20
Batch size	8
Learning rate	3e-5
Gradient clipping	10
α	0.5
β (on SGD)	3
β (on MultiWOZ)	5

Table 1: Hyperparameters used for the experiments in this paper.

for consistency of sampling. We also conducted experiments to verify whether the pre-training to transfer the schema encoding improves the performance and whether SET-DST successfully works on the pre-training step, although the major part in our experiments is the fine-tuning on MultiWOZ including few-shot settings. We evaluated the pre-training performance on KLUE⁴ dataset (Park et al., 2021), which is a Korean dataset for DST, in addition to SGD.

SET-DST needs not only slot information but also a schema. However, MultiWOZ has no schema and no concepts of service and intent; thus, we created a schema for MultiWOZ including services, slots, intents, and corresponding descriptions. Figure 4 shows an example of the schema for MultiWOZ. In our experiments on MultiWOZ, an intent means activated domain. In other words, the system classifies an intent as active when the domain of conversation is changed or a conversation starts. MultiWOZ further has no labels for activated intents, thus we automatically added the labels by detecting the changes of domain.

We further tried to fine-tune the system without intents because it is possible that the concepts of intent are unnatural in MultiWOZ. In this setting, Equation 3, 6, 7, 10, 15 are ignored, I_t is removed from Equation 12, and \mathcal{L}_t^I is removed from Equation 17. In Figure 3a, U_t is replaced with State: { Inform - restaurant_location - San Jose }, and in Figure 3b, U_t is replaced with State:

⁴<https://github.com/KLUE-benchmark/KLUE>.

	JA
TRADE (Wu et al., 2019)	45.60%*
DSTQA (Zhou and Small, 2019)	51.17%
LABES-S2S (Zhang et al., 2020b)	51.45%
DST-Picklist (Zhang et al., 2020a)	53.30%
MinTL-BART (Lin et al., 2020)	53.62%
TripPy (Heck et al., 2020)	55.29%
SimpleTOD (Hosseini-Asl et al., 2020)	55.76%
PPTOD (Su et al., 2021)	57.45%
ConvBERT-DG (Mehri et al., 2020)	58.70%
TripPy+SCoRe (Yu et al., 2020)	60.48%
TripPy+CoCoAug (Li et al., 2020)	60.53%
TripPy+SaCLog (Dai et al., 2021)	60.61%
SET-DST (Ours)	60.39%
SET-DST w/o intent	62.07%

Table 2: DST results on the test set of MultiWOZ in joint accuracy. *: the result is reported by Eric et al. (2020).

```
{ Inform - destination_city - San
Francisco ; Inform - origin_city
- Chicago }.
```

A dataset schema can be variously defined depending on the developer, and our goal is to make the system represent any schema for DST. In the pre-training step, we augmented the schema of SGD dataset to avoid overfitting to the given schema. The schema provides names of services, slots, and intents with short descriptions. We defined some alternatives of the names and descriptions, and sampled inputs for the schema encoder from the augmented schema. Figure 5 shows some examples of the alternatives for bank service.

In dialogue state, multiple slots and intents can be activated at once. However, the order has no meaning in dialogue state. The state generator is trained to generate the dialogue state based on textual label, so it is possible that the order causes wrong optimization and overfitting. Thus, we shuffled the order of slots and intents when making the labels.

We used BERT-base-uncased⁵ model for the schema encoder and GPT-2⁶ model for the state generator. In our experiments, the pre-training step took about two days, and the fine-tuning step took about a day on a TitanRTX GPU. Table 1 lists hyperparameters that are used in our experiments.

⁵<https://huggingface.co/bert-base-uncased>.

⁶<https://huggingface.co/gpt2>.

Few-shot rate		JA	
		w/ intent	w/o intent
w/ pre-training	100%	60.39%	62.07%
	30%	53.43%	56.43%
	25%	53.07%	55.61%
	20%	52.73%	54.37%
	15%	40.41%	51.29%
w/o pre-training	10%	31.20%	29.91%
	100%	58.37%	59.10%
	30%	31.08%	48.96%
	25%	30.39%	30.28%
	20%	22.80%	22.35%
	15%	19.38%	17.09%
	10%	10.21%	15.10%

Table 3: Evaluation results on few-shot settings with considering pre-training.

		JA
SGD	Baseline (Rastogi et al., 2020)	43.40%
	SET-DST (Ours)	55.56%
KLUE	Baseline (Park et al., 2021)	50.22%
	SET-DST (Ours)	57.61%

Table 4: Pre-training results on SGD and KLUE compared to their baselines.

4.3 Experimental Results

Table 2 compares the evaluation results of SET-DST to the previous methods on the test set of MultiWOZ. In our experiments, SET-DST achieved new state-of-the-art joint accuracy when fine-tuned without intent.

Table 3 shows the evaluation results on few-shot settings and the improvement by pre-training. When we used less training data, the pre-training with schema encoding was more effective for DST. SET-DST performed reasonably well with only about 20% of the training data. In most cases, the models fine-tuned without intents achieved higher joint accuracy on MultiWOZ.

Table 4 shows the pre-training results on SGD and KLUE. SET-DST outperformed the baselines. These results demonstrate that SET-DST successfully performs DST with just pre-training. In the experiment on KLUE, we used KLUE-BERT⁷ and KoGPT-2⁸ that are large-scale language models pre-trained on Korean corpus.

⁷<https://huggingface.co/klue/bert-base>.

⁸<https://huggingface.co/skt/kogpt2-base-v2>.

5 Discussion

Schema Encoding Development of training dataset for task-oriented dialogue systems including DST requires a lot of cost, especially in low-resource language. Thus, recently, pre-trained models have been widely used to improve the efficiency of training by transferring the knowledge that is learned on large corpus. In this study, our goal is to transfer a pre-trained DST model to a low-resource domain without limiting the transference as language model level. We pre-trained SET-DST on SGD which is a relatively large dataset and fine-tuned it on MultiWOZ to transfer the schema encoding. As a result, the pre-training significantly improved the accuracy on DST. Pre-trained language models have been already used in many fields. However, our method could tackle general DST beyond language modeling on various domains. We believe that SET-DST can assist the development of DST systems in real world without large dataset on the target domain.

Intent on Fine-tuning In this paper, we define the intents as sub-goals to be achieved through a service. SGD has a schema for dialogues between a virtual assistant and an user. Thus, it is assumed that a system achieves the user’s sub-goals by using APIs, and an intent corresponds to an API. Virtual assistant should tackle various services that could consist of one more intents, e.g., to check account balance and to transfer money in bank service. Unlike that, MultiWOZ has no schema and considers no APIs as intents. Thus, the schema that we temporarily created for experiments in the same form as the schema of SGD could cause confusion in generation of dialogue state. We believe that this is why the results without intent were slightly higher in the experiments. Another reason would be the incorrect labels for intents that we automatically created for experiments on MultiWOZ.

The joint accuracy that has been proposed as an evaluation metric for DST considers only slot-value pairs. However, task-oriented dialogue systems should call APIs of external systems to achieve goals, e.g., to search restaurants and to reserve hotels. The systems that predict only slot-value pairs would be insufficient to replace rule-based traditional systems in real-world. Even though use of intents made no improvement in joint accuracy, we believe that encoding the schema including intents is meaningful in terms of approaching more realis-

tic DST.

6 Conclusion

Transfer learning that makes it possible to apply a pre-trained model to new domains has been attempted a lot. However, the attempts for DST have been just to use large-scale pre-trained models. In this paper, we have proposed SET-DST, which is an effective method for DST with transfer learning by using schema encoding. We have demonstrated how to encode the schema for transferable DST and how to use the schema representation for dialogue state generation. Our experiments show that the schema encoding improves joint accuracy even in few-shot settings.

Even though our approach could perform DST well on target domain with few-shot settings, it required some new data to be fine-tuned. As part of our future work, we plan to design a DST model for zero-shot settings.

References

- Lu Chen, Boer Lv, Chi Wang, Su Zhu, Bowen Tan, and Kai Yu. 2020. Schema-guided multi-domain dialogue state tracking with graph attention neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7521–7528.
- Yinpei Dai, Hangyu Li, Yongbin Li, Jian Sun, Fei Huang, Luo Si, and Xiaodan Zhu. 2021. Preview, attend and review: Schema-aware curriculum learning for multi-domain dialog state tracking. *arXiv preprint arXiv:2106.00291*.
- 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.
- 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 *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 422–428.
- Michael Heck, Carel van Niekerk, Nurul Lubis, Christian Geishauser, Hsien-Chin Lin, Marco Moresi, and Milica Gasic. 2020. Trippy: A triple copy strategy for value independent neural dialog state tracking. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 35–44.

530	Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. <i>arXiv preprint arXiv:2005.00796</i> .	584
531		585
532		586
533		587
534	Sungdong Kim, Sohee Yang, Gyuwan Kim, and Sang-Woo Lee. 2020. Efficient dialogue state tracking by selectively overwriting memory. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 567–582.	588
535		589
536		590
537		591
538		
539	Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	592
540		593
541		594
		595
		596
542	Hwaran Lee, Jinsik Lee, and Tae-Yoon Kim. 2019. Sumbt: Slot-utterance matching for universal and scalable belief tracking. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5478–5483.	597
543		598
544		599
545		600
546		601
547	Shiyang Li, Semih Yavuz, Kazuma Hashimoto, Jia Li, Tong Niu, Nazneen Rajani, Xifeng Yan, Yingbo Zhou, and Caiming Xiong. 2020. Coco: Controllable counterfactuals for evaluating dialogue state trackers. <i>arXiv preprint arXiv:2010.12850</i> .	602
548		603
549		604
550		605
551		606
		607
552	Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, and Pascale Fung. 2020. Mintl: Minimalist transfer learning for task-oriented dialogue systems. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3391–3405.	608
553		609
554		610
555		611
556		612
557		
558	Shikib Mehri, Mihail Eric, and Dilek Hakkani-Tur. 2020. Dialogue: A natural language understanding benchmark for task-oriented dialogue. <i>arXiv preprint arXiv:2009.13570</i> .	613
559		614
560		615
561		616
		617
562	Elnaz Nouri and Ehsan Hosseini-Asl. 2018. Toward scalable neural dialogue state tracking model. <i>arXiv preprint arXiv:1812.00899</i> .	618
563		619
564		
565	Sungjoon Park, Jihyung Moon, Sungdong Kim, Won Ik Cho, Jiyoung Han, Jangwon Park, Chisung Song, Jun-seong Kim, Yongsook Song, Taehwan Oh, et al. 2021. Klue: Korean language understanding evaluation. <i>arXiv preprint arXiv:2105.09680</i> .	620
566		621
567		622
568		623
569		624
		625
570	Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayan-deh, Lars Liden, and Jianfeng Gao. 2021. Soloist: Building task bots at scale with transfer learning and machine teaching. <i>Transactions of the Association for Computational Linguistics</i> , 9:907–824.	626
571		627
572		628
573		629
574		630
		631
575	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	632
576		633
577		634
		635
578	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.	636
579		637
580		638
581		639
582		640
583		
	Liliang Ren, Kaige Xie, Lu Chen, and Kai Yu. 2018. Towards universal dialogue state tracking. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2780–2786.	
	Yixuan Su, Lei Shu, Elman Mansimov, Arshit Gupta, Deng Cai, Yi-An Lai, and Yi Zhang. 2021. Multi-task pre-training for plug-and-play task-oriented dialogue system. <i>arXiv preprint arXiv:2109.14739</i> .	
	Zhuoran Wang and Oliver Lemon. 2013. A simple and generic belief tracking mechanism for the dialog state tracking challenge: On the believability of observed information. In <i>Proceedings of the SIGDIAL 2013 Conference</i> , pages 423–432.	
	Jason D Williams. 2014. Web-style ranking and slu combination for dialog state tracking. In <i>Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)</i> , pages 282–291.	
	Chien-Sheng Wu, Andrea Madotto, Ehsan Hosseini-Asl, Caiming Xiong, Richard Socher, and Pascale Fung. 2019. Transferable multi-domain state generator for task-oriented dialogue systems. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 808–819.	
	Tao Yu, Rui Zhang, Alex Polozov, Christopher Meek, and Ahmed Hassan Awadallah. 2020. Score: Pre-training for context representation in conversational semantic parsing. In <i>International Conference on Learning Representations</i> .	
	Jianguo Zhang, Kazuma Hashimoto, Chien-Sheng Wu, Yao Wang, S Yu Philip, Richard Socher, and Caiming Xiong. 2020a. Find or classify? dual strategy for slot-value predictions on multi-domain dialog state tracking. In <i>Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics</i> , pages 154–167.	
	Yichi Zhang, Zhijian Ou, Min Hu, and Junlan Feng. 2020b. A probabilistic end-to-end task-oriented dialog model with latent belief states towards semi-supervised learning. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9207–9219.	
	Victor Zhong, Caiming Xiong, and Richard Socher. 2018. Global-locally self-attentive encoder for dialogue state tracking. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1458–1467.	
	Li Zhou and Kevin Small. 2019. Multi-domain dialogue state tracking as dynamic knowledge graph enhanced question answering. <i>arXiv preprint arXiv:1911.06192</i> .	
	Su Zhu, Jieyu Li, Lu Chen, and Kai Yu. 2020. Efficient context and schema fusion networks for multi-domain dialogue state tracking. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings</i> , pages 766–781.	