

# Data Augmentation for Copy-Mechanism in Dialogue State Tracking

Xiaohui Song<sup>1</sup>, Liangjun Zang<sup>1</sup>, Yipeng Su<sup>1</sup>, Xing Wu<sup>3</sup>, Jizhong Han<sup>1</sup> and Songlin Hu<sup>1,2</sup>

<sup>1</sup>Institute of Information Engineering, Chinese academy of Sciences, Beijing, China

<sup>2</sup>School of Cyber Security, University of Chinese academy of Sciences, Beijing, China

<sup>3</sup>Baidu Inc., Beijing, China

{songxiaohui, zangliangjun, suyipeng, hanjizhong, husonglin}@iie.ac.cn, wuxing03@baidu.com

## Abstract

While several state-of-the-art approaches to dialogue state tracking (DST) have shown promising performances on several benchmarks, there is still a significant performance gap between seen slot values (*i.e.*, values that occur in both training set and test set) and unseen ones (values that occur in training set but not in test set). Recently, the copy-mechanism has been widely used in DST models to handle unseen slot values, which copies slot values from user utterance directly. In this paper, we aim to find out the factors that influence the generalization ability of a common copy-mechanism model for DST. Our key observations include: 1) the copy-mechanism tends to memorize values rather than infer them from contexts, which is the primary reason for unsatisfactory generalization performance; 2) greater diversity of slot values in the training set increase the performance on unseen values but slightly decrease the performance on seen values. Moreover, we propose a simple but effective algorithm of data augmentation to train copy-mechanism models, which augments the input dataset by copying user utterances and replacing the real slot values with randomly generated strings. Users could use two hyper-parameters to realize a trade-off between the performances on seen values and unseen ones, as well as a trade-off between overall performance and computational cost. Experimental results on three widely used datasets (WoZ 2.0, DSTC2, and Multi-WoZ 2.0) show the effectiveness of our approach.

## 1 Introduction

Task-oriented dialogue system interacts with users in natural language to accomplish tasks such as restaurant reservation or flight booking. The goal of dialogue state tracking is to provide a compact representation of the conversation at each dialogue turn, called *dialogue state*, for the system to decide the next action to take. A typical dialogue state consists of goal of user, action of the current user utterance (*inform*, *request* etc.) and dialogue history [Young *et al.*, 2013]. All of these are defined in a particularly designed *ontology*

that restricts which slots the system can handle, and the range of values each slot can take. Tracking the user’s goal is the focus of this task. To accomplish the tracking task, most DST models take the user’s utterance at the current turn, a slot to track and dialogue history as input, output the corresponding value if the user triggers the input slot. Considering an example of restaurant reservation, users can *inform* the system some restrictions of their goals (*e.g.*, `inform(food = thai)`) or *request* further information they want (*e.g.*, `request(phone number)`) at each turn.

**User:**I’m looking for a restaurant that serves thai food.

`state:inform(food=thai)`

**System:**There are two, one in the west end and one in the centre of town. Do you have preference?

**User:**The one on the west end, please. Can I have the phone number?

`state:inform(food=thai, area=west)`

`request(phone number)`

Having access to an ontology that contains all possible values simplifies the tracking problem in many ways. However, in a real-world dialogue system, it is often impossible to enumerate all possible values for each slot. To reduce the dependence on the ontology, PtrNet[Xu and Hu, 2018] uses the Pointer Network[Vinyals *et al.*, 2015] to handle the unknown values that are not defined in the ontology, since then the attention-based *copy-mechanism* inspired by Pointer Network become widely used in state-of-the-art DST approaches[Ren *et al.*, 2019; Wu *et al.*, 2019; Gao *et al.*, 2019]. The copy-mechanism based DST models directly copy slot values from the dialogue history, thus reducing the need to pre-define all slot values in the ontology.

However, there is still a significant performance gap between seen slot values (*i.e.*, values that occur in both training set and test set) and unseen ones<sup>1</sup> (*i.e.*, values that occur in only training set but not test set). In this paper, our goal is two-fold: 1) we figure out which factors influence the generalization ability of copy-mechanism based models; 2) we develop a simple but effective algorithm to improve the performance of predicting unseen slot values.

<sup>1</sup>In this paper, we use ‘unseen’ values and ‘unknown’ values interchangeably.

We conduct two experiments to figure out the reason for the poor generalization performance of copy-mechanism based models. In the first experiment, we replace all slot values in either test set or training set with randomly generated values, and we observe that F1-scores drop dramatically in each experimental setting. Consequently, we conclude that the model attempts to memorize the slot values rather than infer them from contexts because their contexts keep the same. In the second experiment, to take full advantage of contexts of slot values, we augment a dataset by copying user utterances several times and replacing all slot values with randomly generated strings. We observe that the generalization performance is positively related to the diversity of slot values (*i.e.*, the number of unique values for slots), and that the performance on unseen values is very close to seen values. Consequently, we conclude that the model learns contextual information for slot values and infers correctly them from their contexts.

Based on the above observations, we proposed a simple and effective method of data augmentation to train copy-mechanism models. The algorithm augments the input dataset by copying user utterances and replacing the real slot values with randomly generated strings. Importantly, it could determine the optimal size of a synthetic training set automatically. Experimental results show that the synthetic training sets greatly improve the overall performances of copy-mechanism models. Moreover, users could use two hyper-parameters to realize a trade-off between the performances on seen values and unseen ones, as well as a trade-off between overall performance and computational cost.

The rest of this paper is organized as follows: we first review the recent advances in both DST and the copy mechanism in Section 2. Then, we describe the datasets we used and a general copy-mechanism based model in Section 3. Our data analysis process and key observations are described in Section 4. We propose our data augmentation approach to DST task and evaluate its performance in Section 5. Finally, we conclude our work and discuss the future work.

## 2 Related Works

### 2.1 Dialogue State Tracking

Deep-learning has recently shown its power to the dialogue state tracking challenges [Williams *et al.*, 2013; Henderson *et al.*, 2014a; Henderson *et al.*, 2014b]. Neural Belief Tracker (NBT) [Mrkšić *et al.*, 2016] applied representation learning to learn features as opposed to hand-crafting features. Ptr-Net [Xu and Hu, 2018] aimed to handle unknown values. GLAD [Zhong *et al.*, 2018] addressed the problem of extraction of rare slot-value pairs. The number of parameters of these models increased with the number of slots. [Ramadan *et al.*, 2018] tried to share parameters across slots, but the model had to iterate all slots and values defined in the ontology at each dialogue turn. [Rastogi *et al.*, 2017] generated a fixed set of candidate values using a separate SLU module but suffered from error accumulation. TRADE[Wu *et al.*, 2019] was a simple copy-augmented generative model that tracked dialogue states without ontology and enabled zero-shot and few-shot DST in a new domain. [Gao *et al.*, 2019] used the pretrained language model *BERT*[Devlin *et al.*, 2018] and

copy-mechanism to predict explicitly expressed values.

### 2.2 Copy-Mechanism

Copy-mechanism in deep learning is a general concept, which means an output is copied from an input sequence. This idea was first proposed in Pointer Network[Vinyals *et al.*, 2015]. The Pointer Network is designed to learn the conditional probability of an output sequence with elements that are discrete tokens corresponding to positions in an input sequence. It can address the problems such as sorting variable sized sequences and various combinatorial optimization. Inspired by Pointer Network, CopyNet[Gu *et al.*, 2016] first incorporated copy-mechanism in sequence-to-sequence learning, and the Pointer-Generator Network[See *et al.*, 2017] combined the Pointer Network and a generator to retain the ability to produce novel words. As mentioned in the introduction, Xu and Hu first reformulated DST problem to take advantage of the flexibility enabled by Pointer Network. Then, the copy-mechanism become the most common sub-structure of the DST models, which is used in several state-of-the-art DST models such as TRADE[Wu *et al.*, 2019], COMER[Ren *et al.*, 2019], etc. DSTRead[Gao *et al.*, 2019] used a span prediction method proposed in DrQA[Chen *et al.*, 2017] for machine reading task, which is also an application of the copy mechanism.

## 3 Datasets, Baseline Model and Evaluation

In this section, we first describe the datasets we used, then present a baseline model that implements the copy mechanism widely used in DST models, and finally present evaluating metric.

### 3.1 Datasets

Our datasets are extracted from three datasets widely used in DST tasks, *i.e.*, WoZ 2.0[Wen *et al.*, 2016], DSTC2[Henderson *et al.*, 2014a] and Multi-WoZ 2.0[Budzianowski *et al.*, 2018]. To test the performance of copy mechanism on unseen slot values, we focus on the slots of which the values are non-enumerable. Table 1 lists the slots, the number of slot values, and the number of samples in the following experiments.

datasets	slots	values train(total)	samples train(test)
WoZ	food	73(75)	2536(1646)
DSTC2	food	72(73)	11677(9890)
Multi	hotel-name	35(37)	60027(73720)
	train-destination	19(20)	
	train-departure	23(25)	
	attraction-name	88(95)	
	taxi-destination	198(214)	
	taxi-departure	188(210)	
	restaurant-name	131(139)	
	restaurant-food	94(95)	
	bus-departure	1(1)	
bus-destination	1(1)		

Table 1: Slots we use in experiments in WoZ 2.0, DSTC2 and Multi-WoZ datasets. Samples include negative samples for the slot gate.

Each sample corresponds to one dialogue turn, along with a **slot**, its **active** state, and its corresponding **value**. An example is as follows:

**utterance:**I'm looking for a restaurant that serves thai food.  
**slot:** food, **active:**True, **value:** thai

The active attribute of a slot will be `True` (corresponding to a positive sample) if the user triggers the slot, otherwise it will be set `False` (corresponding to a negative sample). In the test set, each turn of dialogue will be paired with each slot to reach the convincing results.

As shown in Table 1, there are a lot of overlapping values over slots between the training set and test set. To observe the performance on unseen slot values, we design a function of random string generation to generate new unseen slot values. Then, on the basis of the random string generation function, we define another function of generating synthesis datasets.

### Random String Generation

The function `randstr(strlen)` takes string length as input and output a randomly generated string. Algorithm 1 describes its generating process. It randomly samples chars from a fix char sequence, we define the char sequence as 26 lowercase letters and 3 spaces<sup>2</sup>. Spaces are used to generate multi-words values, for example, when `strlen` is 10, the probability of generating multi-words values is about  $1 - (26/29)^{10} \approx 0.66$ .

---

#### Algorithm 1 Generate Random String(`randstr`)

---

**Input:** length of random string  $L$

**Output:** a random string with length  $L$

- 1: define a char sequence  $S = [a - z, ' ', ' ', ' ']$ .
  - 2: **while** True **do**
  - 3:   Generate  $L$  random integers from the discrete uniform distribution as indexes  $idx$ .
  - 4:   Concatenate all chars indexed by  $idx$  in char sequence  $S \rightarrow s$ .
  - 5:   Remove leading and trailing spaces of  $s$ .
  - 6:   **if**  $\text{len}(s) == L$  **then**
  - 7:     break
  - 8:   **end if**
  - 9: **end while**
  - 10: **return**  $s$
- 

### Synthesis Dataset Construction

The synthesis dataset construction function  $DC(D, n, \theta)$  takes an original dataset as input and outputs a new dataset as follows. It first duplicate the data samples in which `active=True` in  $D$  into  $n$  copies, get  $D'$ . For each data sample in  $D'$  that contains a active slot, the function replaces the corresponding value with a random value generated by `randstr`, with a probability of  $\theta$ . Specially, all values in set will be replaced if  $\theta = 1$ , and  $DC(D, 1, 0)$  keep the input dataset unchanged.

<sup>2</sup>in this paper we use 3 spaces and `strlen=10` to generate more natural multi-words values

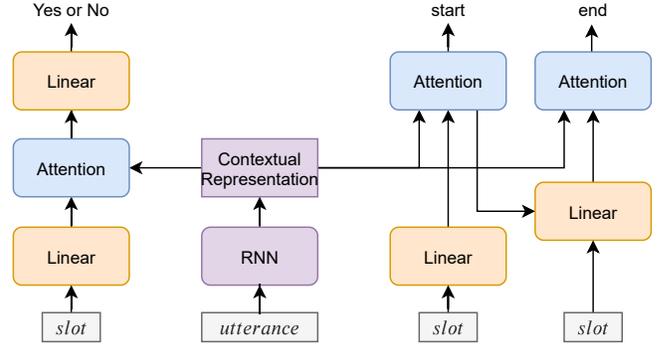


Figure 1: The architecture of our baseline model. The inputs are embedding of slot, word embeddings of user and system utterance ( $W_1, W_2, \dots, W_n$ ). It is a common sub-structure in SOTA DST models.

### 3.2 Baseline Model

Without loss of generality, we implement a basic copy-mechanism based model. The model takes an utterance  $U$  and a slot  $s$  as input, and output whether  $s$  is active and the positions of its corresponding value if  $s$  is active. The model consists of three important components: utterance encoder, attention-based copy-mechanism, and slot gate. Its architecture is presented in Figure 1.

#### Utterance Encoder

For utterance  $U$ , we simply concatenate user and system utterance ( $U^{usr}$  and  $U^{sys}$ ) by a particular symbol `<usr>`.

$$U = U^{sys} \oplus \langle \text{usr} \rangle \oplus U^{usr} \quad (1)$$

We use a bidirectional LSTM [Hochreiter and Schmidhuber, 1997] to get the contextual representation  $H^P$  of  $U$ .

$$H^P = \text{BiLSIM}^P(U_{emb}) \in \mathbb{R}^{n \times d_h} \quad (2)$$

where  $n$  denotes the number of words in  $U$ ,  $U_{emb} \in \mathbb{R}^{n \times d_{emb}}$  the word embeddings of  $U$ ,  $d_{emb}$  the dimension of embeddings, and  $d_h$  the dimension of LSTM hidden states.

#### Attention-based Copy-Mechanism

We define an attention function  $\text{Attn}(Q, V)$  to calculate attention scores between query feature  $Q \in \mathbb{R}^{d_{emb}}$  and context feature  $V \in \mathbb{R}^{n \times d_h}$ . The computing process is as follows: 1) do a linear transform for both  $Q$  and  $V$ , 2) use the dot product result as attention scores, 3) normalize through `softmax` function.

$$Q' = QW_q, V' = VW_v, \alpha_i = Q'V'_i \quad (3)$$

$$\text{scores}_i = \exp \alpha_i / \sum_j \exp \alpha_j \quad (4)$$

$$\text{contexts} = \sum_i \text{scores}_i V'_i \quad (5)$$

We use a single linear layer ( $\text{Linear}(X) = WX + b$ ) to encode slots, and then use the attention function defined above to calculate both start and end positions of its value.

$$s_{enc} = \text{Linear}_{slot}(s_{emb}) \quad (6)$$

$$p_{enc} = \text{Linear}_p(s_{enc}) \quad (7)$$

$$\text{contexts}^p, \text{scores}^p = \text{Attn}_{span}(p_{enc}, H^P) \quad (8)$$

$$q_{enc} = \text{Linear}_q(s_{enc} \oplus \text{contexts}^p) \quad (9)$$

$$\text{contexts}^q, \text{scores}^q = \text{Attn}_{span}(q_{enc}, H^P) \quad (10)$$

$$\text{start} = \arg \max_j \text{scores}_j^p \quad (11)$$

$$\text{end} = \arg \max_j \text{scores}_j^q \quad (12)$$

### Slot Gate

A binary classifier is used to determine whether a slot is triggered by a user or not, where the single linear layer  $\text{Linear}_{cls2}$  produces a probability over  $[0, 1]$  based on attention contexts. A slot is triggered if  $cls_{result} > 0.5$ .

$$s_{cls} = \text{Linear}_{cls1}(s_{emb}) \quad (13)$$

$$\text{contexts}^{cls}, \text{scores}^{cls} = \text{Attn}_{cls}(s_{cls}, H^P) \quad (14)$$

$$\alpha_{cls} = \text{Linear}_{cls2}(\text{contexts}^{cls}) \quad (15)$$

$$cls_{result} = \text{sigmoid}(\alpha_{cls}) \quad (16)$$

### Implementation Details

To enhance the effectiveness of the experiments, all experiments in this paper share the same settings. We use randomly initialized word embeddings of dimension 300 with dropout[Srivastava *et al.*, 2014] rate 0.5. The model is trained with Adam[Kingma and Ba, 2014] optimizer and the learning rate is  $1e-4$ .

In all experiments, we train the baseline model for 80 epochs on WoZ dataset, 32 epochs on DSTC2 dataset and 32 epochs on Multi-WoZ dataset to guarantee good convergence. At each epoch, we evaluate the model on the dev set and save the checkpoint and select the best to get the final result on the test set when the training process completes. All experiments are conducted on a single NVIDIA RTX 2080Ti GPU.

### 3.3 Evaluation Metric

We use F1 scores as the primary metric in all experiments in this paper. Specifically, if a slot is determined to be active by the slot gate and the model predicts the correct value (both the start and end positions) for the slot, then it is a true positive sample. If the slot gate outputs `False`, then the extracted value will be ignored.

## 4 Data Analysis and Observations

In this Section, we analyze the factors that influence the generalization performance of the baseline model.

### 4.1 The model tends to memorize values

As shown in TRADE[Wu *et al.*, 2019], the slots that share similar values or have correlated values have similar embeddings. We guess that the attention based copy-mechanism is inclined to memorize values that appear in the train set rather than infer them from contexts. To verify our argument, we

conduct four groups of experiments using original and synthetic datasets. The experimental results are presented in Table 2, where the synthetic training/test sets are built using the dataset construction function  $DC(D, 1, 1)$  (see definition in section 3.1) under constraint of sharing the same set of randomly generated values.

	dataset	original test	synthetic test
original train	WoZ2.0	0.9241	0.5801
	DSTC2	0.9850	0.3231
	MultiWoZ	0.9079	0.4412
synthetic train	WoZ2.0	0.5283	0.9586
	DSTC2	0.2385	0.9863
	Multi-WoZ	0.3631	0.8837

Table 2: F1 scores on three datasets. The synthetic training/test sets are constructed using  $DC(D, 1, 1)$  with the same set of slot values.

As shown in Table 2, the model perform well when both training and test set are original or synthetic, but perform poorly when one is original and the other is synthetic. The former corresponds to the case of seen slot values, and the latter corresponds to the case of unseen slot values. Since we do not change the contexts of slot values, we conclude that the model attempts to memorize the slot values rather than infer them from contexts. Interestingly, the model shows pretty good performances in the case that training/test set are both synthetic, which confirms that it is reasonable to use random strings to augment datasets.

### 4.2 Greater diversity of values brings better generalization

We have argued that the model tends to memorize the values that appear in the train set. In other words, the model pay little attention to context information of slot values, which is a very useful feature for the model to recognize the correct value for a slot. Hence, we would like to make full use of context information of slot values. Intuitively, high diversity of slot values make the model more difficult to learn from slot values, and duplicated contexts make the model easier to learn from less varied contexts. Consequently, an interesting question is:

*If we greatly increase the diversity of slot values and duplicate contexts in the train set, will the model prefer inferring slot values from slot contexts?*

### Greater diversity of values improves generalization

We could adjust the diversity of slot values by controlling the number of slot values in the training set. Let  $N$  be the number of all slot values (may duplicate) appeared in the training set. We set the numbers of unique values with  $[20, 20 \times \alpha^1, \dots, 20 \times \alpha^{29} = N]$  (*i.e.*, totally 30 numbers), where 20 and  $N$  is the minimum and maximum number of unique slot values respectively. We sample in this way to present a clear trend for small number of values. For each

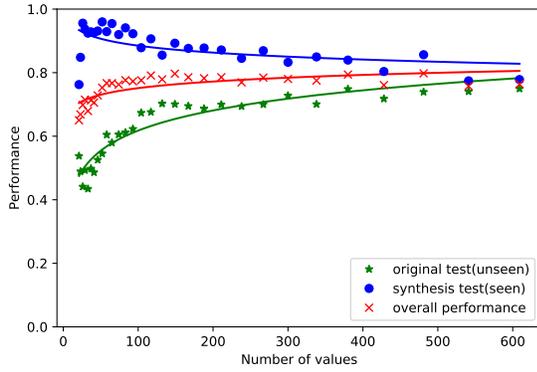


Figure 2: Performance on dataset constructed using  $DC(WoZ, 1, 1)$  with different size of values set, values in training and modified test set are all randomly generated from the same values set, *all values in original test set are unseen values*. Overall performance in the figure is the mean of F1 score at the two test sets. Experiments also conducted on DSTC2 and Multi-WoZ dataset and got similar results.

number, we first use the function of random string generation (see algorithm 1) to produce a new set of slot values, and then use the function  $DC(D, n = 1, \theta = 1)$  of dataset construction to produce a synthetic training set and a synthetic test set with the same set of slot values. In this case, all values in the synthetic test set could be found in the synthetic training set, and all values in the original test set are completely unseen for the synthetic training set.

As shown in Figure 2, the F1-scores of original test set (*i.e.*, unseen values) and average F1-scores increase rapidly with increasing number of unique slot values. Hence, we conclude that the diversity of values are positively correlated to model’s generalization ability.

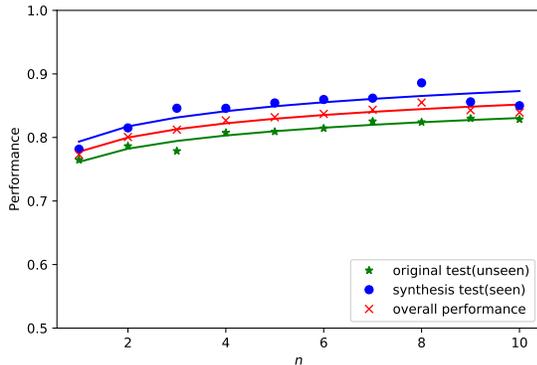


Figure 3: Performance on dataset constructed by function  $DC(WoZ, n, 1)$ ,  $n \in [1..10]$ . **Note** that in this figure, we didn’t control the size of values set, the size of values set increased as  $n$  increased.

### More samples further improves generalization

In Figure 2, the number of different values is less than or equal to the number of samples that contains active slots. What will happen if we continue to increase the number of different values? In this experiment, we use the function  $DC(D, n, \theta = 1)$  to produce synthetic training/test set,

where  $n$  is taken from  $[1..10]$  and all generated values are different.

As shown in Figure 3, the F1-scores of original/synthetic test sets and average F1-scores increase with increasing number of unique slot values. Hence, we conclude that the number  $n$  of copies are positively correlated to model’s generalization ability.

To sum up, we have similar conclusion from Figure 2 and Figure 3. That is, with increasing number of slot values, the performance on unseen values and the overall performance increase substantially but the increasing speed gets slower. Therefore, increasing the diversity of values of the training set to some extent is a feasible way to improve the model’s generalization. In addition, we can find that the performance on seen values (*i.e.*, synthetic test set) is always better than that on unseen values (*i.e.*, original test set).

## 5 Data Augmentation for DST

In this section, we first propose a simple data augmentation algorithm based on the conclusions in Section 4, and then evaluate our algorithm experimentally.

---

### Algorithm 2 Data Augmentation Algorithm

---

**Input:** datasets  $D, T_s, T_u$ , randomly initialized tracker  $T$   
**Hyperparameters:** the proportion of the new value  $\theta$ ,  $\text{eps}=\epsilon$   
**Output:** the times  $n$  that copies original dataset, the tracker  $T$  trained on  $DC(D, n, \theta)$ , the performance  $P$  of  $T$ .

- 1: Let the number of copies  $t = 1$ , performance  $res=[]$
  - 2: train  $T$  on  $D$ , and get an overall performance  $P$  on  $T_s$  and  $T_u$  follow steps in section 3.2
  - 3:  $res \leftarrow P$
  - 4: **while** True **do**
  - 5: re-initialize the tracker  $T$
  - 6:  $D' = DC(D, n, \theta)$
  - 7: train  $T$  on  $D'$ , get the overall performance  $P$  on  $T_s$  and  $T_u$  the same as step 2
  - 8:  $res \leftarrow P$
  - 9:  $n = n * 2$
  - 10: **if**  $\text{len}(res) < 3$  **then**
  - 11: continue
  - 12: **end if**
  - 13: **if**  $res[-1] - res[-2] < \epsilon$  and  $res[-2] - res[-3] < \epsilon$  **then**
  - 14: break
  - 15: **end if**
  - 16: **end while**
  - 17: **return**  $n, T, \max(res)$
- 

### 5.1 Approach

Given a training set  $D$ , a test set  $T_s$  containing all seen values and another test set  $T_u$  containing all unseen values, our algorithm outputs the optimal number  $n$  for the optimization problem as follows:

$$\arg \min_{n \geq 2} \{P_\theta(2n) - P_\theta(n) \leq \epsilon \text{ and } P_\theta(n) - P_\theta(n/2) \leq \epsilon\} \quad (17)$$

where  $P_\theta$  is the average of the performances on  $T_s$  and  $T_u$ .

There are two hyper-parameters  $\theta$  and  $\epsilon$  in our algorithm. Firstly,  $\theta$  is used to control the proportion of unseen values when we generate a training set using the function  $DC(D, n, \theta)$ . Users can use  $\theta$  to trade-off the performances between seen values and unseen ones. For example, a user may set  $\theta$  with 0.5 when he believe that seen values and unseen values are equally important. Secondly,  $\epsilon$  is used to control the terminating condition of the algorithm. Even though the overall performance is positively related to the number  $n$  of copies in  $DC(D, n, \theta)$ , the computational cost also increases when the size of training set increases. Users can use  $\epsilon$  to trade-off between performance and computational cost.

**This problem can be defined to find a point in Figure 3 that has the best overall performance.** We use the doubling method to search the best  $n$  to construct the training set within an acceptable computational cost. Algorithm 2 presents the details of this process.

## 5.2 Performance

We evaluate our data augmentation algorithm on three datasets: WoZ 2.0, DSTC2 and Multi-WoZ. We focus on slots that have non-enumerable values (see Table 1) and present experimental results in Table 3. Here we use  $\theta = 0.5$  to balance the ratio between seen and unseen values, and set the precision  $\epsilon$  to 0.01 to get best overall performance within limited computational resources.

From Table 3 we can find that, with our data augmentation algorithm, the performance of the baseline model on unseen values improves significantly, and the performance on seen values decrease slightly.

Dataset	Model	seen	unseen
WoZ 2.0	baseline	0.9241	0.5801
	+DA(n=64)	0.9073(↓1.7%)	0.8818(↑30.17%)
DSTC2	baseline	0.9850	0.3231
	+DA(n=32)	0.9711(↓1.39%)	0.9591(↑63.6%)
Multi-WoZ	baseline	0.9079	0.4412
	+DA(n=32)	0.9032(↓0.65%)	0.8847(↑44.35%)

Table 3: F1 scores of our data augmentation(DA) approach on three datasets’ test set(seen) and modified test set(unseen), with  $\theta = 0.5$  and  $\epsilon = 0.01$ . And  $n$  is the number of copies in the dataset construct function  $DC(D, n, \theta)$  that gains the best overall performance in Algorithm 2.

## 5.3 Hyperparametric Analysis

We introduce two hyperparameters in our data augmentation algorithm,  $\theta$  and  $\epsilon$ .  $\theta$  is used in the dataset construct function  $DC(D, n, \theta)$  to control the proportion of new generated values in the train set, and  $\epsilon$  is used in the terminate condition of the doubling searching to trade off the performance and computational cost. The relation between performance and  $\theta$  is shown in Figure 4, and the effect of  $\epsilon$  is shown in Table 4.

In Figure 4, we apply different  $\theta$  to our algorithm, and record the performance  $P_s$  and  $P_u$  for each  $\theta$ . We can find that **a larger  $\theta$  increases the performance on unseen values but decreases the performance on seen values**. As we described in section 4.2, performance on seen values always

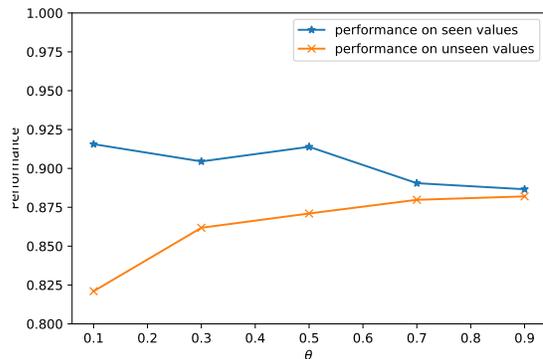


Figure 4: The effect of  $\theta$  on WoZ dataset with  $\epsilon = 0.01$ .

better than that on unseen values. This figure confirm that the proportion of new generated values could be used to balance the performance on seen and unseen values.

As the number of search steps increases, the performance gain gets smaller, which can be concluded from Figure 3. From table 4, we could observe that **a small  $\epsilon$  leads to more searching steps and better performance, while a large  $\epsilon$  means lower computational cost and slightly worse performance**. Thus, we could use  $\epsilon$  to realize a trade-off between performance and computational cost.

$\epsilon$	search steps	overall performance
0.01	7(t=64)	0.8946
0.02	6(t=32)	0.8871
0.03	6(t=32)	0.8871
0.04	4(t=8)	0.8443
0.05	3(t=4)	0.8443

Table 4: The effect of the  $\epsilon$  on the WoZ dataset with  $\theta = 0.5$ .

## 6 Conclusion and Future Work

This paper focuses on the problem of improving generalization ability of copy-mechanism based models for DST task, especially for the slots that have non-enumerable values. Our conclusions include: Firstly, the copy-mechanism model attempts to memorize values rather than infer them from contexts, which is the crucial reason for unsatisfactory generalization performance. Secondly, the model’s generalization improves dramatically with increasing diversity of slot values. Thirdly, data augmentation for copy-mechanism models is feasible and effective in improving generalization of these models.

In future work, we can use the character-level features of the slot values since the values for the same slot often share somehow similar spelling. For example, an address or location often contains uppercase letters, and a slot about time often consists of Arabic numerals and symbol ‘:’. In addition, the effect of contexts’ diversity remains unexplored, which may help reduce computational cost further.

## References

- [Budzianowski *et al.*, 2018] Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. Multiwoz—a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. *arXiv preprint arXiv:1810.00278*, 2018.
- [Chen *et al.*, 2017] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, 2017.
- [Devlin *et al.*, 2018] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [Gao *et al.*, 2019] Shuyang Gao, Abhishek Sethi, Sanchit Agarwal, Tagyoung Chung, and Dilek Hakkani-Tur. Dialog state tracking: A neural reading comprehension approach. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, pages 264–273, 2019.
- [Gu *et al.*, 2016] Jiatao Gu, Zhengdong Lu, Hang Li, and Victor OK Li. Incorporating copying mechanism in sequence-to-sequence learning. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1631–1640, 2016.
- [Henderson *et al.*, 2014a] Matthew Henderson, Blaise Thomson, and Jason D Williams. 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, 2014.
- [Henderson *et al.*, 2014b] Matthew Henderson, Blaise Thomson, and Jason D Williams. The third dialog state tracking challenge. In *2014 IEEE Spoken Language Technology Workshop (SLT)*, pages 324–329. IEEE, 2014.
- [Hochreiter and Schmidhuber, 1997] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [Kingma and Ba, 2014] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [Mrkšić *et al.*, 2016] Nikola Mrkšić, Diarmuid O Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve Young. Neural belief tracker: Data-driven dialogue state tracking. *arXiv preprint arXiv:1606.03777*, 2016.
- [Ramadan *et al.*, 2018] Osman Ramadan, Paweł Budzianowski, and Milica Gašić. Large-scale multi-domain belief tracking with knowledge sharing. *arXiv preprint arXiv:1807.06517*, 2018.
- [Rastogi *et al.*, 2017] Abhinav Rastogi, Dilek Hakkani-Tür, and Larry Heck. Scalable multi-domain dialogue state tracking. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 561–568. IEEE, 2017.
- [Ren *et al.*, 2019] Liliang Ren, Jianmo Ni, and Julian McAuley. Scalable and accurate dialogue state tracking via hierarchical sequence generation. 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 1876–1885, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [See *et al.*, 2017] Abigail See, Peter J Liu, and Christopher D Manning. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368*, 2017.
- [Srivastava *et al.*, 2014] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [Vinyals *et al.*, 2015] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in Neural Information Processing Systems*, pages 2692–2700, 2015.
- [Wen *et al.*, 2016] Tsung-Hsien Wen, David Vandyke, Nikola Mrksic, Milica Gasic, Lina M Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. A network-based end-to-end trainable task-oriented dialogue system. *arXiv preprint arXiv:1604.04562*, 2016.
- [Williams *et al.*, 2013] Jason Williams, Antoine Raux, Deepak Ramachandran, and Alan Black. The dialog state tracking challenge. In *Proceedings of the SIGDIAL 2013 Conference*, pages 404–413, 2013.
- [Wu *et al.*, 2019] Chien-Sheng Wu, Andrea Madotto, Ehsan Hosseini-Asl, Caiming Xiong, Richard Socher, and Pascale Fung. Transferable multi-domain state generator for task-oriented dialogue systems. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2019.
- [Xu and Hu, 2018] Puyang Xu and Qi Hu. An end-to-end approach for handling unknown slot values in dialogue state tracking. *arXiv preprint arXiv:1805.01555*, 2018.
- [Young *et al.*, 2013] Steve Young, Milica Gašić, Blaise Thomson, and Jason D Williams. Pomdp-based statistical spoken dialog systems: A review. *Proceedings of the IEEE*, 101(5):1160–1179, 2013.
- [Zhong *et al.*, 2018] Victor Zhong, Caiming Xiong, and Richard Socher. Global-locally self-attentive encoder for dialogue state tracking. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1458–1467, 2018.