

# TripPy: A Triple Copy Strategy for Value Independent Neural Dialog State Tracking

Michael Heck, Carel van Niekerk, Nurul Lubis,  
Christian Geishauser, Hsien-Chin Lin, Marco Moresi, Milica Gašić

Heinrich Heine University Düsseldorf, Germany

{heckmi, niekerk, lubis, geishaus, linh, moresi, gasic}@hhu.de

## Abstract

Task-oriented dialog systems rely on dialog state tracking (DST) to monitor the user’s goal during the course of an interaction. Multi-domain and open-vocabulary settings complicate the task considerably and demand scalable solutions. In this paper we present a new approach to DST which makes use of various copy mechanisms to fill slots with values. **Our model has no need to maintain a list of candidate values. Instead, all values are extracted from the dialog context on-the-fly. A slot is filled by one of three copy mechanisms:** (1) Span prediction may extract values directly from the user input; (2) a value may be copied from a system inform memory that keeps track of the system’s inform operations; (3) a value may be copied over from a different slot that is already contained in the dialog state to resolve coreferences within and across domains. Our approach combines the advantages of span-based slot filling methods with memory methods to avoid the use of value picklists altogether. We argue that our strategy simplifies the DST task while at the same time achieving state of the art performance on various popular evaluation sets including Multiwoz 2.1, where we achieve a joint goal accuracy beyond 55%.

## 1 Introduction

The increasing popularity of natural language human-computer interaction urges the development of robust and scalable task-oriented dialog systems. In order to fulfill a user goal, a dialogue system must be capable of extracting meaning and intent from the user input, and be able to keep and update this information over the continuation of the dialog (Young et al., 2010). This task is called dialog state tracking (DST). Because the next dialog system action depends on the current state of the conversation, accurate dialog state tracking (DST) is absolutely vital.

U: i'm looking for an **expensive** restaurant in the **center** of town.  
S: there are 33 restaurants [...]. would you like to narrow your search by type of food?  
U: i **do not care** about food. surprise me.  
S: **fitzbillies restaurant** serves british food, [...]. would you like to hear about any others? [...].  
U: **that sounds fine**, can you book it please and get me the reference number?  
S: sure, what day and time would you like and how many people?  
U: i would like a table for **5** at **11:30** on **tuesday** [...]  
S: okay, the booking was successful, [...]. is there anything else i can help you with?  
U: i'm also looking for a place to stay. it needs [...]**free wifi** and [be] in the **same area as the restaurant**.

| Turn | Domain-slot pair      | Value       | Type                    | Coreference     |
|------|-----------------------|-------------|-------------------------|-----------------|
| 0    | restaurant-pricerange | expensive   | span                    |                 |
| 0    | restaurant-area       | center      | span                    |                 |
| 1    | restaurant-food       | <dontcare>  | (dontcare)              |                 |
| 2    | restaurant-name       | fitzbillies | informed                |                 |
| 3    | restaurant-people     | 5           | span                    |                 |
| 3    | restaurant-book_time  | 11:30       | span                    |                 |
| 3    | restaurant-book_day   | tuesday     | span                    |                 |
| 4    | hotel-internet        | <true>      | (bool)                  |                 |
| 5    | hotel-area            | center      | coreference (multiturn) | restaurant-area |

Figure 1: Example dialog in MultiWOZ.

DST is tasked to extract from the user input information on different concepts that are necessary to complete the task at hand. For example, in order to recommend a restaurant to a user, the system needs to know their preferences in terms of price, location, etc. These concepts are encapsulated in an ontology, where dialogue domain (e.g., restaurant or hotel), slot (e.g., price range or location), and value (e.g. cheap or expensive) are defined. Solving this information extraction task is prerequisite for forming a belief over the dialog state.

Traditional approaches to DST operate on a fixed ontology and perform prediction over a pre-defined set of slot-value pairs (Mrkšić et al., 2016; Liu and Lane, 2017; Zhong et al., 2018). Such approaches perform very well on datasets which are defined over fairly small ontologies. Apply these methods to more complex datasets however reveals various limitations (Ren et al., 2018; Nouri and Hosseini-Asl, 2018). First, it is often difficult to obtain a complete ontology for a task. Second, slot-value pairs that were outside the ontology or the training data are impossible to capture during test time.

Third, such methods at best scale linearly with the size of the ontology. **Most importantly, the idea of fixed ontologies is not sustainable, as in real world applications they are subject to constant change.**

Human-computer interactions often need to be defined over multiple domains at the same time, ideally with unrestricted vocabulary. Recent approaches to multi-domain and open-vocabulary DST extract values from the dialog context directly by predicting value spans in the input (Gao et al., 2019; Chao and Lane, 2019; Kim et al., 2019; Zhang et al., 2019). Span prediction is a demonstrably potent method to detect relevant information in utterances, but its major drawback is that it only suits *extractive* values that are explicitly expressed as a sequence of tokens. This is the reason why span-based methods benefit from the support of a picklist, i.e., a list of value candidates from which a system can choose. Still, these methods fall short when handling nuanced and subtle phenomena that often occur in natural conversations, such as coreference and value sharing (“I’d like a hotel in the same area as the restaurant.”), and implicit choice (“Any of those is ok.”).

In this work, we propose a new approach to value independent multi-domain DST:

1. In addition to extracting values directly from the user utterance via span prediction and copy, our model creates and maintains two memories on-the-fly, one for system inform slots, and one for the previously seen slots.
2. The *system inform memory* solves the implicit choice issue by allowing copy mechanism from concepts mentioned by the system, e.g., values that are offered and recommended.
3. The *DS memory* allows the use of values already existing in the dialogue state to infer new values, which solves the coreference and value sharing problems.

We call this approach **TripPy**, **Triple copy** strategy DST.<sup>1</sup> Our experiments results show that our model is able to handle out-of-vocabulary and rare values very well during test time, demonstrating good generalization. In a detailed analysis we take a closer look at each of the model’s components to study their particular roles.

<sup>1</sup>Our code will be released upon publication of this work.

## 2 Related Work

Dialog state tracking has been of broad interest to the dialog research community, which is reflected by the existence of a series of DST challenges (Henderson et al., 2014; Rastogi et al., 2019). These challenges consistently pushed the boundaries of DST performance. Current state-of-the-art has to prove to work on long, diverse conversations in multiple domains with a high slot count and principally unrestricted vocabulary (Eric et al., 2019). Dialogs of such complex nature are tough for traditional approaches that rely on the availability of a candidate list due to scalability and generalization issues (Mrkšić et al., 2016; Liu and Lane, 2017; Ramadan et al., 2018; Rastogi et al., 2017).

Span-based approaches recently alleviated both problems to some extent. Here, slot values are extracted from the input directly by predicting start and end positions in the course of the dialog. For instance, Xu and Hu (2018) utilizes an attention-based recurrent network with a pointer mechanism to extract values from the context. This extractive approach has its limitations, since many expressible values are not found verbatim in the input, but rather mentioned implicitly, or expressed by a variety of rephrasings.

With the assistance of contextual models such as BERT (Devlin et al., 2018), issues arising from expressional variations can be mitigated. Recent work has demonstrated that encoding the dialog context with contextual representations supports span prediction to generalize over rephrasings. SUMBT (Lee et al., 2019) utilizes BERT to encode slot IDs and candidate values and learns slot-value relationships appearing in dialogs via an attention mechanism. Dialog context is encoded with recurrence. BERT-DST (Chao and Lane, 2019) employs contextual representations to encode each dialog turn and feeds them into classification heads for value prediction. The dialog history, however, is not considered for slot filling. In Gao et al. (2019), DST is rendered as a reading comprehension task that is approached with a BERT-based dialog context encoder. A slot carryover prediction model determines whether previously detected values should be kept in the DS for the current turn.

An alternative to span prediction is value generation. TRADE (Wu et al., 2019) and MA-DST (Kumar et al., 2020) generate a DS from the input using a copy mechanism to combine the distributions over a pre-defined vocabulary and the vocabulary

of current context. SOM-DST (Kim et al., 2019) applies a similar mechanism for value generation, but takes the previous dialog turn as well as the previous DS as input to BERT to predict the current DS. A state operation predictor determines whether a slot actually needs to be updated or not. The downside of generative models is that they tend to produce invalid values, for instance by word repetitions or omissions.

Recently, a hybrid approach called DS-DST has been proposed that makes use of both span-based and picklist-based prediction for slot-filling (Zhang et al., 2019). In contrast to generative approaches, picklist-based and span-based methods use existing word sequences to fill slots. DS-DST somewhat alleviates the limitations of span prediction by filling a subset of slots with a picklist method instead.

Recent works seemed to reveal a trade-off between the level value independence in a model and the DST performance. Chao and Lane (2019) and Gao et al. (2019) solely rely on span-prediction, but their performance lacks behind methods that at least partially rely on a pre-defined list of candidate values. This has impressively been demonstrated by Zhang et al. (2019). Their model could not compete when relying on span-prediction entirely. In contrast, when relying solely on their picklist slot-filling method, they achieved the to-date best performance on MultiWOZ 2.1. The proposed dual-strategy approach lies favorably between these two extremes.

To the best of our knowledge, none of the recent approaches to complex DST tasks such as MultiWOZ (Budzianowski et al., 2018; Eric et al., 2019) are value independent in the strict sense. What’s more, they tremendously benefit from the use of a value candidate list. Our work tackles this limitation by introducing a triple copy strategy that relies on span-prediction as well as memory mechanisms. In contrast to other hybrid approaches such as Zhang et al. (2019), our memory mechanisms create candidate lists of values on-the-fly with the dialog context as only source of information, thus avoiding the use of pre-defined picklists. We let the model decide which strategy to choose for each slot at each turn. Our approach differs from Chao and Lane (2019) and Kim et al. (2019) in that we consider the dialog history as context in addition to the current turn. We also differ from approaches like Lee et al. (2019) since we do not employ recurrence. Like Kim et al. (2019), we use auxiliary

inputs at each turn, but we do so as a late feature fusion strategy. With our slot-value copy mechanism to resolve coreferring value phrases, we employ a method which is reminiscent of Gao et al. (2019)’s slot carryover, but with the sharp distinction that we copy values between different slots, facilitating value sharing within and across domains.

### 3 TripPy: Triple Copy Strategy for DST

Our model expects the following input format to perform dialog state tracking. Let  $X = \{(U_1, M_1), \dots, (U_T, M_T)\}$  be the sequence of turns that comprise a dialog of length  $T$ .  $U_t$  is the user utterance at turn  $t$ ,  $M_t$  is the system utterance that precedes the user utterance. The task of the model is (1) to determine for every turn whether any of the  $N$  domain-slot pairs in  $S = \{S_1, \dots, S_N\}$  is present, (2) to predict the values for each  $S_n$  and (3) to track the dialog state  $DS_t$  over the course of the dialog, i.e., for  $t \in [1, T]$ .

We employ a triple-copy strategy to fill the slots. The intuition is that values are either explicitly expressed by the user, that they are expressed by the system and referred to by the user via confirmation or rejection, or that they have been expressed earlier in the dialog as assignment to another domain-slot pair (coreference). Each of these cases is handled by one of three copy mechanisms. It becomes apparent that slots can not be filled by exclusively resorting to one particular copy method. Therefore, we employ slot gates that determine at each turn which method to use to fill the respective slot.

Figure 2 depicts our model. We encode the dialog context with a BERT front-end and feed-forward the resulting contextual representations to various classification heads to solve the sub-tasks for DST. The aggregate sequence representation is the input to the slot gates. The sequence of token representations is the input to the span predictors.

#### 3.1 Context Encoder

We use BERT (Devlin et al., 2018) as front-end to encode at each turn  $t$  the dialog context as

$$R_t = \text{BERT}([\text{CLS}] \oplus U_t \oplus [\text{SEP}] \oplus M_t \oplus [\text{SEP}] \oplus H_t \oplus [\text{SEP}]), \quad (1)$$

where  $H_t = (U_{t-1}, M_{t-1}), \dots, (U_1, M_1)$  is the history of the dialog up to and excluding turn  $t$ . The special token [CLS] precedes every input sequence to BERT, and [SEP] separates portions of the input sequence. It is then  $R_t = [r_t^{\text{CLS}}, r_t^1, \dots, r_t^{\text{seq}_{\max}}]$ ,

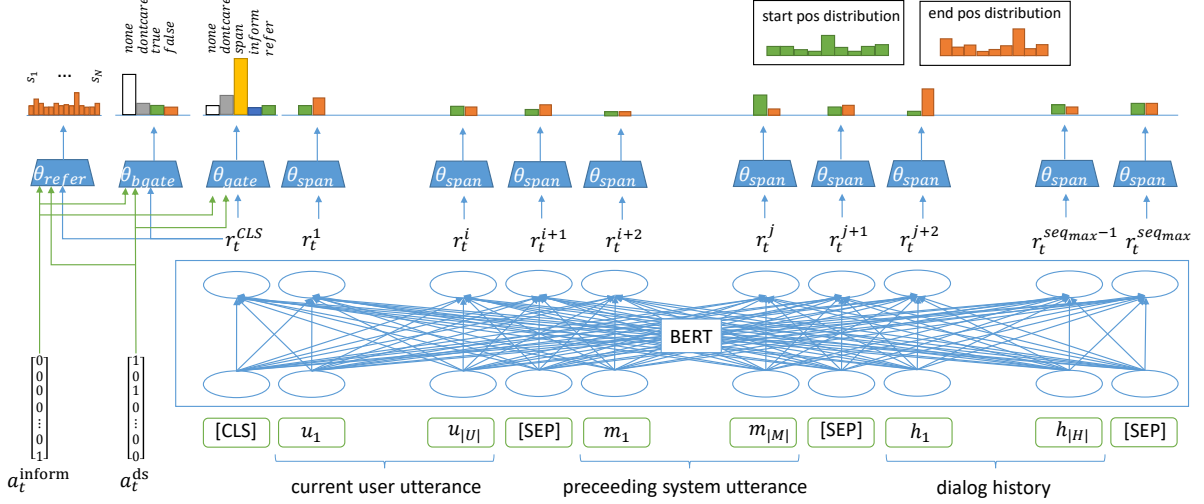


Figure 2: Architecture of our proposed model. TripPy takes the turn and dialog history as input and outputs a DS.

where  $r_t^{\text{CLS}}$  is a representation of the entire turn including the dialog context  $H_t$ . The vectors  $r_t^1$  to  $r_t^{\text{seqmax}}$  are contextual representations for the sequence of input tokens (including special tokens). Both types of representations are used for the following classification tasks.

### 3.2 Slot Gates

Our model is equipped with a slot gate for each domain-slot pair. This ensures greatest flexibility for multi-domain DST, as there is no restriction as to how many domains might be present in a single turn. At each turn  $t$ , slot gates assign each slot  $S_n$  to one of the classes in  $C = \{none, dontcare, span, inform, refer\}$ . The first two labels express special cases. *none* denotes that the slot does not take a value in this turn and *dontcare* states that any value is acceptable for this slot. The remaining three labels each denote one of the model’s copy mechanisms. *span* indicates that a value is present in  $U_t$  that can be extracted via span prediction. *inform* indicates that the user refers to a value that has been uttered by the system in  $M_t$ . Lastly, *refer* indicates that the user refers to a value that is already present in  $DS_t$ .

The input to the slot gates is  $r_t^{\text{CLS}}$ , and the probability distribution over classes  $C$  for domain-slot pair  $S_n$  at turn  $t$  is  $p_{t,s}^{\text{gate}}(r_t^{\text{CLS}}) =$

$$\text{softmax}(W_s^{\text{gate}} \cdot r_t^{\text{CLS}} + b_s^{\text{gate}}) \in \mathbb{R}^5, \quad (2)$$

i.e., each slot gate is realized by a trainable linear layer classification head for BERT.

Boolean slots, i.e., slots that only take binary values, are treated separately. Here, the list of possible

classes is  $C_{\text{bool}} = \{none, dontcare, true, false\}$  and the slot gate probability is  $p_{t,s}^{\text{bgate}}(r_t^{\text{CLS}}) =$

$$\text{softmax}(W_s^{\text{bgate}} \cdot r_t^{\text{CLS}} + b_s^{\text{bgate}}) \in \mathbb{R}^4. \quad (3)$$

### 3.3 Span-based Value Prediction

For each slot  $s$  that is to be filled via span prediction, a domain-slot specific span prediction layer takes the token representations  $[r_t^1, \dots, r_t^{\text{seqmax}}]$  of the entire dialog context for turn  $t$  as input and projects them as follows:

$$[\alpha_{t,i}^s, \beta_{t,i}^s] = W_s^{\text{span}} \cdot r_t^i + b_s^{\text{span}} \in \mathbb{R}^2 \quad (4a)$$

$$p_{t,s}^{\text{start}} = \text{softmax}(\alpha_{t,i}^s) \quad (4b)$$

$$p_{t,s}^{\text{end}} = \text{softmax}(\beta_{t,i}^s) \quad (4c)$$

$$\text{start}_t^s = \text{argmax}(p_{t,s}^{\text{start}}) \quad (4d)$$

$$\text{end}_t^s = \text{argmax}(p_{t,s}^{\text{end}}). \quad (4e)$$

Each span predictor is realized by a trainable linear layer classification head for BERT, followed by two parallel softmax layers to predict start and end position. Note that there is no special handling for erroneously predicting  $\text{end}_t^s < \text{start}_t^s$ . In practice, the resulting span will simply be empty.

### 3.4 System Inform Memory for Value Prediction

The system inform memory  $I_t = \{I_t^1, \dots, I_t^N\}$  keeps track of all slot values that were informed by the system in dialog turn  $t$ . A slot in  $DS_t$  needs to be filled by an informed value, if the user positively refers to it, but does not express the value such that span prediction can be used. E.g., in Figure 1 the



slot gate for domain-slot `<restaurant, name>` should predict *inform*. The slot is filled by copying the informed value into the dialog state, i.e.,  $DS_t^s = I_t^s$ , where  $i$  is the index of the respective domain-slot.

### 3.5 DS Memory for Coreference Resolution

The more complex a dialog can be, the more likely it is that coreferences need to be resolved. For instance, the name of a restaurant might very well be the destination of a taxi ride, but the restaurant might not be referred to explicitly upon ordering a taxi within the same conversation. Coreference resolution is challenging due to the rich variety of how to form referrals, as well as due to the fact that coreferences often span multiple turns. An example of a coreference that can be handled by our model is found in the example in Figure 1.

The third copy mechanism utilizes the DS as a memory to resolve coreferences. If a slot gate predicts that the user refers to a value that has already been assigned to a different slot during the conversation, then the probability distribution over all possible slots that can be referenced is

$$p_{t,s}^{\text{refer}}(r_t^{\text{CLS}}) = \text{softmax}(W_{\text{refer}}^s \cdot r_t^{\text{CLS}} + b_{\text{refer}}^s) \in \mathbb{R}^{N+1}, \quad (5)$$

i.e., for each slot, a linear layer classification head either predicts the slot which contains the referenced value, or *none* for no reference.

### 3.6 Auxiliary Features

Some recent approaches to neural DST utilize auxiliary input to preserve contextual information. For instance, SOM-DST adds the dialog state to its single-turn input as a means to preserve context across turns.

We already include contextual information in the input to BERT by appending the dialog history  $H_t$ . In addition to that, we also create auxiliary features based on the system inform memory and the DS memory. We generate two binary vectors  $a_t^{\text{inform}} \in \{0, 1\}^N$  and  $a_t^{\text{ds}} \in \{0, 1\}^N$  that indicate whether (1) a slot has recently been informed (based on the system inform memory), or (2) a slot has already been filled during the course of the dialog (based on the DS memory). These vectors are added to the output of BERT in a late fusion approach, and the slot gate probabilities in Equations 2, 3 and 5 become  $p_{t,s}^{\text{gate}}(\hat{r}_t^{\text{CLS}})$ ,  $p_{t,s}^{\text{bgate}}(\hat{r}_t^{\text{CLS}})$  and  $p_{t,s}^{\text{refer}}(\hat{r}_t^{\text{CLS}})$ , with  $\hat{r}_t^{\text{CLS}} = r_t^{\text{CLS}} \oplus a_t^{\text{inform}} \oplus a_t^{\text{ds}}$ .

### 3.7 Partial Masking

We partially mask the dialog history  $H_t$  by replacing values with BERT’s generic [UNK] token. The masking is partial in the sense that it is applied only to the past system utterances. For the system utterances, the contained values are known and their masking is straightforward. The idea behind partially masking the history is that the model is compelled to focus on the historical context information rather than the sighting of specific values. This should result in more robust representations  $r_t^{\text{CLS}}$  and therefore better overall slot gate performance.

### 3.8 Dialog State Update

We employ the same rule-based update mechanism as [Chao and Lane \(2019\)](#) to track the dialog state across turns. At every turn, we update a slot, if a value has been detected which is not *none*. If a slot-value is predicted as *none*, then the slot will not be updated.

## 4 Experimental Setup

### 4.1 Datasets

We train and test our model on four datasets, MultiWOZ 2.1 ([Eric et al., 2019](#)), WOZ 2.0 ([Wen et al., 2016](#)), sim-M and sim-R ([Shah et al., 2018](#)). Among these, MultiWOZ 2.1 is by far the most challenging dataset. It is comprised of over 10000 multi-domain dialogs defined over a fairly large ontology. There are 5 domains (train, restaurant, hotel, taxi, attraction) with 30 domain-slot pairs that appear in all portions of the data.

The other datasets are single-domain and significantly smaller. Evaluations on these mainly serve as sanity check to show that we don’t overfit to a particular problem. Some slots in sim-M and sim-R show a high out-of-vocabulary rate, making them particularly interesting for evaluating value independent DST.

The single domain datasets come with span labels. However, MultiWOZ 2.1 does not. We therefore generate our own span labels by matching the ground truth value labels to their respective utterances.

### 4.2 Evaluation

We compute the joint goal accuracy (JGA) on all test sets for straightforward comparison with other approaches. The joint goal accuracy defined over a dataset is the ratio of dialog turns in that dataset

| Models      | MultiWOZ 2.1       |
|-------------|--------------------|
| DST-reader  | 36.40%             |
| DST-span    | 40.39%             |
| SUMBT       | 42.40%**           |
| TRADE       | 45.60%             |
| MA-DST      | 51.04%             |
| DS-DST      | 51.21%             |
| SOM-DST     | 52.57%             |
| DS-picklist | 53.30%             |
| TripPy      | <b>55.29±0.28%</b> |

Table 1: DST Results on MultiWOZ 2.1 in JGA ( $\pm$  denotes the standard deviation. \*\* MultiWOZ 2.0 result.

for which all slots have been filled with the correct value according to the ground truth. Note that *none* needs to be predicted if a slot value is not present in a turn. In addition to JGA, we compute the accuracy of the slot gates (joint and per-class) and various other metrics for a more detailed analysis of model design decisions.

We run each test three times with different seeds and report the average numbers for more reliable results. MultiWOZ 2.1 is in parts labeled inconsistently. For a fair evaluation, we consider a value prediction correct, if it matches any of its valid labels (for instance "centre" and "center" for the slot-value *hotel-area=centre*) as being correct. We semi-automatically analyzed value label inconsistencies in the training portion of the dataset in order to identify all label variants for any given value. During testing, these mappings are applied as is.

### 4.3 Training

We use the pre-trained *BERT-base-uncased* transformer (Vaswani et al., 2017) as context encoder front-end. This model has 12 hidden layers with 768 units and 12 self-attention heads each. The maximum input sequence length is set to 180 tokens after WordPiece tokenization (Wu et al., 2016), except for MultiWOZ 2.1, where we set this parameter to 512. We compute the joint loss as

$$\mathcal{L} = 0.8 \cdot \mathcal{L}_{\text{gate}} + 0.1 \cdot \mathcal{L}_{\text{span}} + 0.1 \cdot \mathcal{L}_{\text{refer}}. \quad (6)$$

The function for all losses is joint cross entropy. As there is no coreferencing in the evaluated single-domain datasets, the refer loss is not computed in those cases and the loss function is

$$\mathcal{L} = 0.8 \cdot \mathcal{L}_{\text{gate}} + 0.2 \cdot \mathcal{L}_{\text{span}} \quad (7)$$

instead. Span predictors are presented only spans from the user utterances  $U_i$  to learn from (includ-

| Models   | WOZ 2.0          |
|----------|------------------|
| NBT      | 84.2%            |
| BERT-DST | 87.7%            |
| GLAD     | 88.1%            |
| GCE      | 88.5%            |
| StateNet | 88.9%            |
| SUMBT    | 91.0%            |
| TripPy   | <b>92.7±0.2%</b> |

Table 2: DST Results on WOZ 2.0.

| Models   | sim-M            | sim-R            |
|----------|------------------|------------------|
| SMD-DST  | 96.8%*           | 94.4%*           |
| LU-DST   | 50.4%            | 87.1%            |
| BERT-DST | 80.1%            | 89.6%            |
| TripPy   | <b>83.5±1.2%</b> | <b>90.0±0.2%</b> |

Table 3: DST Results on sim-M and sim-R. \* should be considered as oracle because the value candidates are ground truth labels.

ing the user utterances in the history portion  $H_i$  of the input). During training we set the span prediction loss to zero for all slots that are not labeled as *span*. Likewise, the coreference prediction losses are set to zero if slots are not labeled as *refer*. For optimization we use Adam optimizer (Kingma and Ba, 2014) and backpropagate through the entire network including BERT, which constitutes a fine-tuning of the latter. The initial learning rate is set to  $2e^{-5}$ . We conduct training with a warmup proportion of 10% and let the LR decay linearly after the warmup phase. Early stopping is employed based on the JGA of the development set. During training we use dropout (Srivastava et al., 2014) on the BERT output with a rate of 30%. We do not use slot value dropout (Xu and Sarikaya, 2014) except for one dataset (sim-M), where performance was greatly affected by this measure (see Section 5.1).

## 5 Experimental Results

Tables 1, 3 and 2 show the performance of our model in comparison to various baselines. TripPy achieves state-of-the-art performance on all four evaluated datasets, with varying distance to the runner-up. Most notably, we were able to push the performance on MultiWOZ 2.1, the most complex task, by another 2.0% absolute compared to the previous top scoring method, achieving 55.3% JGA. The improvements on the much smaller datasets WOZ 2.0, sim-M and sim-R demonstrate that the model benefits from its design on single-domain

| Model                              | JGA    |
|------------------------------------|--------|
| Span prediction only (entire turn) | 42.63% |
| + triple copy mechanism            | 49.23% |
| + dialog history                   | 52.58% |
| + auxiliary features               | 54.08% |
| + masking                          | 54.29% |
| TripPy (full sequence width)       | 55.29% |

Table 4: Ablation experiments for our model.

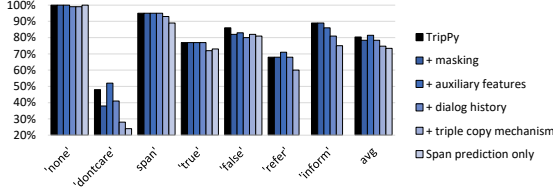


Figure 3: Per class performance of the slot gates for different versions of our model (ablation study).

tasks as well. The following analysis serves a better understanding of our model’s strengths.

### 5.1 Analysis

We analyse the performance of TripPy on ablation experiments on MultiWOZ 2.1 (see Table 4). Our baseline model is best compared to BERT-DST (Chao and Lane, 2019); we only take single turns as input, and only use span prediction to extract values from the turn. The resulting performance is comparable to other span-based methods such as DST-reader and DST-span and confirms that the dialogs in MultiWOZ are too complex to only be handled by this information extracting mechanism alone.

**Impact of the triple copy mechanism** Using our proposed triple copy mechanism pushes the performance close to 50%, surpassing TRADE and closing in on the leading hybrid approaches. Especially the performance of the slot gates benefits from this change (see Figure 3). When looking at the F1 score for the individual classes, one can see that the *span* class benefits from the distinction. It is important to point out that none of the coreferences that our model handles can be resolved by span-prediction alone. This means that otherwise guaranteed misses can now be avoided and coreferences can be resolved by copying values between slots. What’s more, using the dialog state memory to resolve coreferences helps value detection across multiple turns, as a value that has been referred to in the current turn might have been assigned to

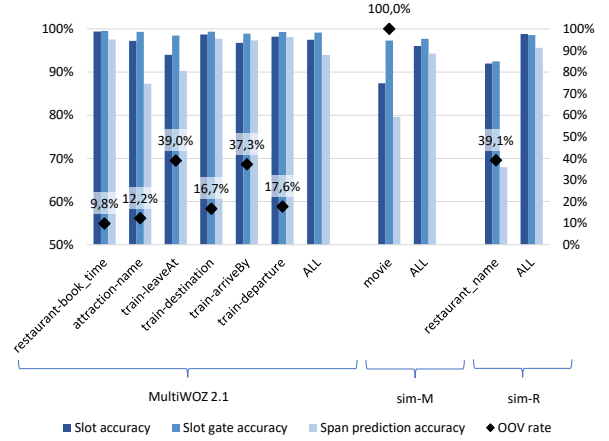


Figure 4: Performance of TripPy on slots with high OOV rate. *ALL* denotes the average of all slots of the respective dataset.

another slot multiple turns before.

**Impact of the dialog history** We found that using the dialog history as additional context information is critical to a good performance, as it reduces contextual ambiguity. This is clearly reflected in the improved performance of the slot gates (see Figure 3, which has two positive effects. First, the presence and type of values is recognized correctly more often. Especially the special value *dontcare*, and boolean slots (taking values *true* and *false*) benefit from the additional context. This is only logical, since they are predicted by the slot gate using the representation vector of the [CLS] token. Second, values are assigned to the correct slot more often than without the additional contextual information. With the additional dialog history, we outperform DS-DST and match SOM-DST, which set the previous state-of-the-art.

**Impact of the auxiliary features** SOM-DST uses single turns as input, but preserves additional contextual information throughout the dialog by using the dialog state as auxiliary input. By adding our memory based auxiliary features in a late fusion approach, we surpass SOM-DST, and ultimately DS-picklist, which performs slot-filling with the knowledge of the full ontology. Even though our features carry less information, that is, only the identities of the informed slots – tracked by the system inform memory – and the identities of the previously seen slots – tracked by the DS memory –, we see substantial improvement using them. Obviously, more information about the progress of the dialog helps the slot gates and the referral gates in

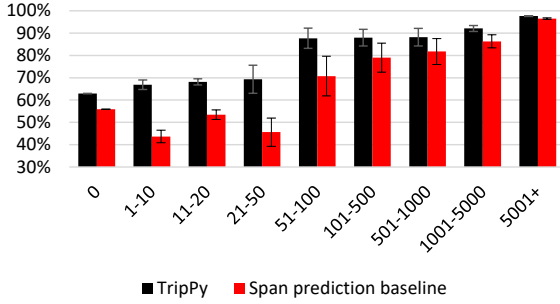


Figure 5: Recall of values depending on the amount of samples seen during training. 0 seen samples means the value is OOV during test time.

their classification tasks.

**Impact of partial masking** We found that masking the informed values in past system utterances does not give a clear benefit, but it also does not harm performance of the slot gates. While the *inform* cases are detected more accurately, some other cases suffer from the loss of information in the input. Overall, there is a minor overall improvement observable. We report the numbers for MultiWOZ in Table 4 and Figure 3, but would like to note that we have seen the same trend on all other datasets as well.

**Impact of the context width** Our best model utilizes the full width of BERT (512 tokens). This is a clear advantage for longer dialogs. Maximal context width is not a decisive factor for the single-domain datasets, since their dialogs tend to be shorter. As expected, we have not seen any change in performance on these. For MultiWOZ, we gain 1% absolute by maximizing the history length to preserve as much of the dialog history as possible, achieving 55.3% JGA.

## 5.2 Generalization Study

It is important that a DST model generalizes well to previously unseen values. We looked at the performance of our model on slots with exceptionally high out-of-vocabulary rates, of which we identified 8 across the evaluated datasets. Figure 4 plots performance measures for these slots and compares them to the average performance for all slots in the respective datasets. Generally, the slots that expect named entities as values show the lowest accuracy. However, the below-average performance of these slots does not seem to be caused by a particularly high OOV rate. Even at 100%, the *movie* slot of sim-M still performs comparably well. Other slots

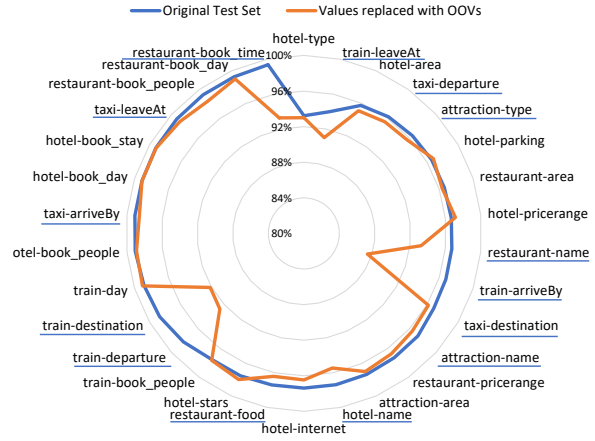


Figure 6: Per-slot accuracy of TripPy on the original test set and the OOV test set. Underlined slot names indicate slots with at least one OOV value.

with relatively high OOV rate still perform close to or better than the average.

Figure 5 plots the recall of values depending on the number of samples seen during training. To our surprise, it does not seem to matter whether a particular value has never been seen during training in order to be detected correctly. OOV values are detected just as well as generally less common values. Our observations however indicate that the model benefits tremendously by seeing a certain minimal amount of training samples for each value, which is somewhere around 50. In other words, if such amounts of data are available, then the model is able to effectively utilize them. In the same Figure we compare TripPy to the span prediction baseline. The latter clearly struggles with OOVs and rare values and generally seems to require more training samples to achieve a good recall. The higher recall on OOV values is likely caused by the fact that many unseen values are of the category time of day, which mostly follows a strict format and is therefore easier to spot. Overall, TripPy clearly generalizes better over sample counts.

To test the limits of our model’s generalization capacities, we manually replaced most of the values in the MultiWOZ test set by (fictional but still meaningful) OOV values. Of the over 1000 unique slot-value pairs appearing in the modified test set, about 84% are OOV after the replacement. Figure 6 compares the per-slot accuracy of our model on the original test set and the OOV test set. Underlined slot names indicate slots with at least one OOV value. Their average OOV rate is 90%. Surprisingly, most of these slots maintain their high



accuracy and only few suffer from the high OOV count. Mainly it is one particular domain, *train*, which suffers above-average performance drops. However, the remainder of the slots maintain their performance. This demonstrates that our model is well equipped to handle OOV values, regardless of the type (e.g., named entity, time of day).

## 6 Conclusion

We have demonstrated that our approach can handle challenging DST scenarios. Having to detect unseen values does not considerably impair our model’s general performance. The information extraction capabilities of our proposed model are rooted in the memory-based copy mechanisms and perform well even in extreme cases as discussed in Section 5.2. The copy mechanisms are not limited by a predefined vocabulary, since the memories themselves are value agnostic.

To further improve the DST capabilities of TripPy, we hope to introduce slot independence as at present its tracking abilities are limited to slots that are predefined in the ontology. For that, We would like to expand our approach towards the schema-guided paradigm for dialog modeling. We also would like to employ a more sophisticated update strategy, for example by adding the option to partially forget. There already exists an intriguing set of works focusing on these issues and we hope to incorporate and expand upon it in the future.

## Acknowledgments

M. Heck, C. van Niekerk and N. Lubis are supported by funding provided by the Alexander von Humboldt Foundation in the framework of the Sofja Kovalevskaja Award endowed by the Federal Ministry of Education and Research, while C. Geishauser, H-C. Lin and M. Moresi are supported by funds from the European Research Council (ERC) provided under the Horizon 2020 research and innovation programme (Grant agreement No. STG2018\_804636).

## References

- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Inigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. *arXiv preprint arXiv:1810.00278*.
- Guan-Lin Chao and Ian Lane. 2019. BERT-DST: Scalable end-to-end dialogue state tracking with bidirectional encoder representations from transformer. *arXiv preprint arXiv:1907.03040*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, and Dilek Hakkani-Tür. 2019. MultiWOZ 2.1: Multi-domain dialogue state corrections and state tracking baselines. *arXiv preprint arXiv:1907.01669*.
- Shuyang Gao, Abhishek Sethi, Sanchit Aggarwal, Tagyoung Chung, and Dilek Hakkani-Tür. 2019. Dialog state tracking: A neural reading comprehension approach. *arXiv preprint arXiv:1908.01946*.
- 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.
- Sungdong Kim, Sohee Yang, Gyuwan Kim, and Sang-Woo Lee. 2019. Efficient dialogue state tracking by selectively overwriting memory. *arXiv preprint arXiv:1911.03906*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Adarsh Kumar, Peter Ku, Anuj Kumar Goyal, Angeliki Metallinou, and Dilek Hakkani-Tür. 2020. MA-DST: Multi-attention based scalable dialog state tracking. *arXiv preprint arXiv:2002.08898*.
- Hwaran Lee, Jinsik Lee, and Tae-Yoon Kim. 2019. SUMBT: Slot-utterance matching for universal and scalable belief tracking. *arXiv preprint arXiv:1907.07421*.
- Bing Liu and Ian Lane. 2017. An end-to-end trainable neural network model with belief tracking for task-oriented dialog. *arXiv preprint arXiv:1708.05956*.
- Nikola Mrkšić, Diarmuid O Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve Young. 2016. Neural belief tracker: Data-driven dialogue state tracking. *arXiv preprint arXiv:1606.03777*.
- Elnaz Nouri and Ehsan Hosseini-Asl. 2018. Toward scalable neural dialogue state tracking model. *arXiv preprint arXiv:1812.00899*.
- Osman Ramadan, Paweł Budzianowski, and Milica Gašić. 2018. Large-scale multi-domain belief tracking with knowledge sharing. *arXiv preprint arXiv:1807.06517*.

- Abhinav Rastogi, Dilek Hakkani-Tür, and Larry Heck. 2017. Scalable multi-domain dialogue state tracking. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 561–568. IEEE.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2019. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. *arXiv preprint arXiv:1909.05855*.
- Liliang Ren, Kaige Xie, Lu Chen, and Kai Yu. 2018. Towards universal dialogue state tracking. *arXiv preprint arXiv:1810.09587*.
- Pararth Shah, Dilek Hakkani-Tür, Gokhan Tür, Abhinav Rastogi, Ankur Bapna, Neha Nayak, and Larry Heck. 2018. Building a conversational agent overnight with dialogue self-play. *arXiv preprint arXiv:1801.04871*.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Milica Gašić, Lina M Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. 2016. A network-based end-to-end trainable task-oriented dialogue system. *arXiv preprint arXiv:1604.04562*.
- 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. *arXiv preprint arXiv:1905.08743*.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Puyang Xu and Qi Hu. 2018. An end-to-end approach for handling unknown slot values in dialogue state tracking. *arXiv preprint arXiv:1805.01555*.
- Puyang Xu and Ruhi Sarikaya. 2014. Targeted feature dropout for robust slot filling in natural language understanding. In *Fifteenth Annual Conference of the International Speech Communication Association*.
- Steve Young, Milica Gašić, Simon Keizer, François Mairesse, Jost Schatzmann, Blaise Thomson, and Kai Yu. 2010. The hidden information state model: A practical framework for POMDP-based spoken dialogue management. *Computer Speech & Language*, 24(2):150–174.
- Jian-Guo Zhang, Kazuma Hashimoto, Chien-Sheng Wu, Yao Wan, Philip S Yu, Richard Socher, and Caiming Xiong. 2019. Find or classify? dual strategy for slot-value predictions on multi-domain dialog state tracking. *arXiv preprint arXiv:1910.03544*.
- Victor Zhong, Caiming Xiong, and Richard Socher. 2018. Global-locally self-attentive dialogue state tracker. *arXiv preprint arXiv:1805.09655*.