

PromptAttack: Probing Dialogue State Tracker with Adversarial Prompts

Xiangjue Dong¹, Yun He^{1*}, Ziwei Zhu^{2*}, James Caverlee¹

¹ Texas A&M University, ² George Mason University

{xj.dong, yunhe, caverlee}@tamu.edu, zzhu20@gmu.edu

Abstract

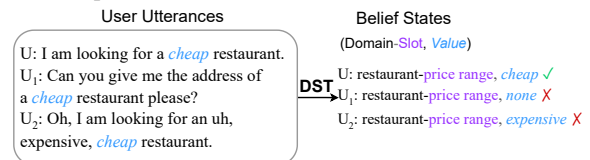
A key component of modern conversational systems is the Dialogue State Tracker (or DST), which models a user’s goals and needs. Toward building more robust and reliable DSTs, we introduce a prompt-based learning approach to automatically generate effective adversarial examples to probe DST models. Two key characteristics of this approach are: (i) it only needs the output of the DST with no need for model parameters, and (ii) it can learn to generate natural language utterances that can target any DST. Through experiments over state-of-the-art DSTs, the proposed framework leads to the greatest reduction in accuracy and the best attack success rate while maintaining good fluency and a low perturbation ratio. We also show how much the generated adversarial examples can bolster a DST through adversarial training. These results indicate the strength of prompt-based attacks on DSTs and leave open avenues for continued refinement.

1 Introduction

Task-oriented dialogue systems aim to help users with tasks through a natural language conversation. Example tasks include booking a hotel or completing a do-it-yourself project. A key component for enabling a high-quality task-oriented dialogue system is the *Dialogue State Tracker* (or DST) which plays an important role in understanding users’ goals and needs (Wu et al., 2019; Hosseini-Asl et al., 2020; Li et al., 2021b; Dai et al., 2021; Feng et al., 2021; Zhao et al., 2021; Balaraman et al., 2021). For example in Figure 1a, given the user utterance “I am looking for a cheap restaurant in the center of the city”, the DST extracts the user’s preference for booking a restaurant, which is typically represented as slot-value pairs such as (restaurant-price range, cheap) and (restaurant-area, center). The current state



(a) Dialogue state tracking task. R represents system response, and U represents user utterance.



(b) Adversarial examples. U₁ and U₂ are adversarial examples based on U which maintain ground-truth values but lead DST models to wrong predictions.

Figure 1: Dialogue examples and adversarial examples.

of the conversation is a primary driver of the subsequent dialogue components (e.g., what is the next action to take? what is the appropriate response to generate?).

For a conversational system designer, it is critical that a deployed DST be robust and reliable, even in the presence of a wide variety of user utterances. Many of these systems are trained over previous user utterances and so may have only limited coverage of the space of these utterances. Further, beyond these benign users, there is also a long history of spammers, trolls, and malicious users who aim to intentionally undermine deployed systems.

Indeed, recent work has demonstrated that careful construction of adversarial examples can cause failures in the DST (Li et al., 2021b; Liu et al., 2021a), leading to incorrect slot-value pairs and degraded user experience. These approaches, however, are mainly hand-crafted or based on heuristics. As a result, there is a research gap in learning-based methods for probing DSTs centered around three key questions: (i) How can we systematically learn effective adversarial examples? (ii) What impact do such discovered examples have on the quality of state-of-the-art DSTs? and (iii) Can we

* Equal Contribution

build more robust DSTs even in the presence of such adversarial examples? Further compounding these questions are the inherent challenges of adversarial examples in the context of a DST: that is, the examples should preserve the semantics of a non-adversarial input while leading to an incorrect prediction *even in the presence of the correct slot-value in the adversarial input* as illustrated in Figure 1b. For example, an adversarial example based on the user utterance “I am looking for a cheap restaurant” that maps to the slot-value pair (restaurant-price range, cheap) should preserve the user intent for “cheap” while leading to the incorrect prediction (restaurant-price range, expensive).

Hence, in this paper, we propose a novel prompt-based learning approach called *PromptAttack* to automatically generate effective adversarial examples to probe DST models. Our approach builds on recent advances in prompt learning, which has demonstrated a strong ability in probing knowledge in pre-trained language models for many NLP tasks (Gao et al., 2021; Li and Liang, 2021; Liu et al., 2021b; Zhu et al., 2022). Concretely, we first show how to find effective adversarial prompts in both a discrete and a continuous setting. In both cases, our approach needs only the output of the DST (e.g., (restaurant-price range, cheap)) with no need for model parameters or other model details. Second, we use the adversarial prompts to generate adversarial examples via a mask-and-filling protocol, resulting in natural language utterances that can be targeted at any DST. As a result, such a prompt-based attack can be widely applied.

Through experiments over four state-of-the-art DSTs and versus competitive baselines, we find that the prompt-based framework leads to the greatest reduction in accuracy for all DSTs, ranging from a 9.3 to 31.0 loss of accuracy of the DST making a correct slot-value prediction. Further, we observe that PromptAttack results in the best attack success rate (that is, how many of the adversarial examples lead to incorrect predictions). Moreover, the generated adversarial examples maintain good fluency and low perturbation ratio, evidence that they are close to legitimate non-adversarial user inputs. We also show how such a prompt-based attack can be used to bolster a DST by augmenting the original training data with adversarial examples, leading to a significant increase in accuracy (from 61.3 to 67.3). These and other results indi-

cate the strength of prompt-based attacks on DSTs and leave open avenues for continued refinement.¹

2 Related Work

Adversarial examples have been widely explored to investigate the robustness of models (Goodfellow et al., 2015). Recent work in the NLP domain has targeted tasks like text classification and inference (Pruthi et al., 2019; Ren et al., 2019; Morris et al., 2020; Jin et al., 2020; Li et al., 2020; Yang et al., 2022; Lei et al., 2022), reading comprehension (Jia and Liang, 2017; Bartolo et al., 2021), named entity recognition (Simoncini and Spanakis, 2021), and machine translation (Belinkov and Bisk, 2018). These works typically aim to construct examples that are imperceptible to human judges while misleading the underlying model to make an incorrect prediction, while also maintaining good fluency and semantic consistency with original inputs (Li et al., 2020). Only a few works have begun to explore adversarial examples in DSTs like CoCo (Li et al., 2021b), which aims to test the robustness of models by creating novel and realistic conversation scenarios. They show that DST models are susceptible to both unseen slot values generated from in and out of the slot domain. Liu et al. (2021a) propose a model-agnostic toolkit to test the robustness of task-oriented dialogue systems in terms of three aspects: speech characteristics, language variety, and noise perturbation. The adversarial examples are based on heuristics and it is unclear how to adapt such an approach to new victim models effectively without more hand-crafted templates. In contrast, we explore in this paper the potential of a learning-based approach to generate effective adversarial examples.

Prompt learning is a recently proposed paradigm for using prompts to better probe and adapt large pre-trained language models (PLMs) to a variety of NLP tasks, e.g., text classification and inference (Gao et al., 2021; Yang et al., 2022; Wang et al., 2022), factual probing (Zhong et al., 2021), summarization (Li and Liang, 2021), and dialogue systems (Madotto et al., 2021; Lee et al., 2021; Zhu et al., 2022; Yang et al., 2023). With the increase in the size of PLMs, prompt learning has been shown to be parameter-efficient (Liu et al., 2021b; He et al., 2022; Lu et al., 2023). There are two types of prompts: discrete (or hard) prompts and continu-

¹Our code is publicly available at <https://github.com/dongxiangjue/PromptAttack>.

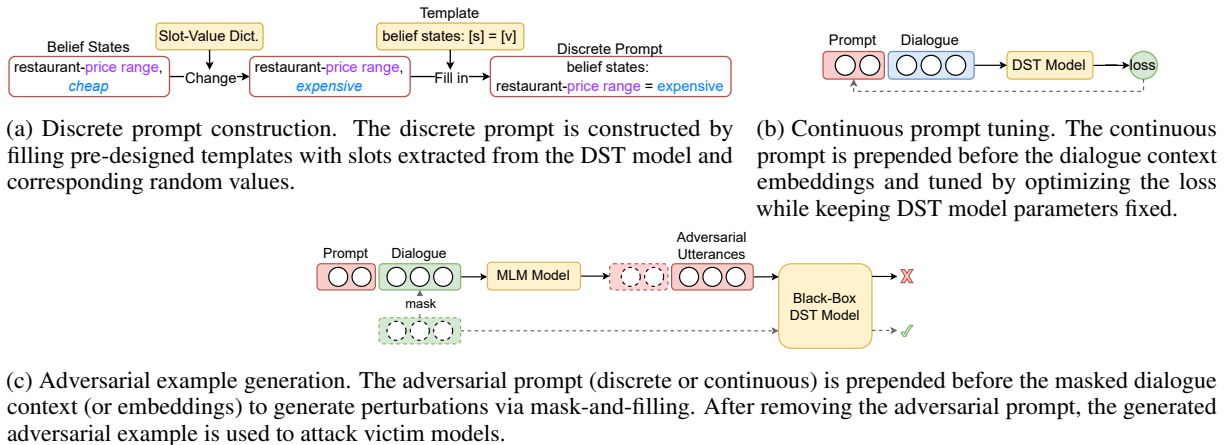


Figure 2: Overview of PromptAttack.

ous (or soft) prompts. Discrete prompts are human-designed text strings (Brown et al., 2020) while continuous prompts are continuous embeddings. Soft prompts proposed by Lester et al. (2021) prepend a sequence of continuous vectors to the input, freeze the language model parameters, and then back-propagate the error during tuning. In this paper, we explore both approaches in the design of our prompt-based attack framework.

Recent works have begun to explore how prompts can be helpful in exposing fundamental flaws in large language models. Yang et al. (2022) shows how to manually design prompts to flip the output of a model for classification tasks. However, it is time-consuming to design and find prompts that are most effective to generate adversarial examples capable of successfully attacking victim models. It is an open question how to leverage prompts for uncovering effective adversarial prompts.

3 PromptAttack

Our prompt-based learning approach proceeds in two stages. First, our goal is to identify adversarial prompts that can effectively probe a DST to reveal gaps in its robustness. In the second, we use these prompts to create adversarial examples that can attack DSTs successfully while maintaining good fluency. Figure 2 shows an overview of the proposed approach. In the following, we first formalize DSTs and the problem of probing a DST. Then, we introduce the details of PromptAttack.

3.1 Task Formulation

DST Task. Let $C_T = \{(r_1, u_1), \dots, (r_T, u_T)\}$ represent a T -turn dialogue, where r_i and u_i ($1 \leq i \leq T$) are the system response and user ut-

terance at the i -th turn, respectively. Each turn (r_i, u_i) contains several slots (e.g., arrive by, leave at) in a specific domain (e.g., taxi), where we denote the N domain-slot pairs as $S = \{s_1, \dots, s_N\}$. At turn t , we denote current user utterance u_t and previous dialogue context $C_t = \{(r_1, u_1), \dots, (r_{t-1}, u_{t-1}), r_t\}$. A DST model aims to extract the dialogue belief state $B_t = \{(s_1, v_1), \dots, (s_N, v_N)\}$ for u_t , where v_j is the associated value for each slot s_j ($1 \leq j \leq N$). For example, given a dialogue (“... I am looking for expensive Mediterranean food.”), the DST model aims to extract expensive for the slot restaurant-price range and Mediterranean for the slot restaurant-food.

Attacking a DST. Given dialogue history C_t , current user utterance u_t , and dialogue belief states B_t , the purpose of an adversarial attack on a DST is to intentionally perturb the original user utterance u_t to get an adversarial example u'_t with the two following characteristics: (i) it should mislead the DST model f to incorrectly predict B'_t , and (ii) it should be fluent in grammar and consistent with the semantics of the original utterance u_t by keeping the slot-value-related information in u_t unchanged. If the adversary can achieve $f(u'_t) = B'_t$, we say the adversarial example u'_t attacks f successfully.

3.2 Finding Adversarial Prompts

We begin by focusing on the first stage of PromptAttack: how to find the most effective adversarial prompts. We explore both discrete prompts (as illustrated in Figure 2a) and continuous prompts (as illustrated in Figure 2b). A discrete prompt approach is a human-designed natural language prompt that is easy to interpret. We pair this with

a treatment of continuous prompts that have more representation capacity.

Discrete Prompt Construction. To begin with, how can we design discrete prompts? For the DST task, it is time-consuming to manually design sentences containing values that are opposite to the ground truth values for each slot as adversarial prompts. Thus, we apply an intuitive template derived from belief states as an adversarial prompt template: “belief states: [s] = [v];”. First, we use the DST model to extract value v_i for each slot s_i in u_t . If v_i is not empty, the corresponding slot name s_i is filled in [s]. Then we pick a random value v'_i from a predefined in-domain Slot-Value Dictionary (Li et al., 2021b) where v'_i and v_i are under the same slot s_i . The new random value v'_i is used to fill the [v] in the template. Thus, the adversarial prompt becomes “belief states: $s_i = v'_i$ ”. As in Figure 2a, given u_t (“I am looking for cheap food.”), the predicted B_t is $\{(restaurant-price\ range, cheap)\}$, then the adversarial prompt is “belief states: restaurant-price range = expensive”, where “expensive” is a random value that is different from the predicted value “cheap”.

Such a template does not have access to true slot-value pairs of the test set and only utilizes the predictions from the victim models. Since the discrete prompts are human-designed, they are more human-readable and easier to interpret. However, to obtain a prompt for each input, victim models must be queried multiple times, which may be unrealistic in some scenarios. Hence, we take the next step to search for better prompts in the embedding space of the model. Specifically, we directly optimize the continuous input embedding space through continuous prompt tuning to find the adversarial prompt vectors that are most effective.

Continuous Prompt Tuning. Continuous prompts are input-agnostic sequences of embeddings with tunable parameters that are optimized directly in the continuous embedding space of the model, as shown in Figure 2b. In our task, the length of continuous prompt \mathbf{p}_{att} is m , denoted as $\mathbf{p}_{att} = \mathbf{p}_1 \dots \mathbf{p}_m$ where each $\mathbf{p}_i \in \mathbb{R}^d (1 \leq i \leq m)$ is a dense vector with the same dimension d as the DST’s input embedding (e.g., 768 for TripPy). Given the initialization of \mathbf{p}_{att} , we concatenate it with the representation of user utterance \mathbf{e}_u and update it by keeping all other model parameters fixed and optimize the loss of the training set. To find the adversarial prompts \mathbf{p}_{att} that could lead DST

models f to wrong predictions B'_t effectively, we maximize the loss for the ground truth belief states B_t for all user utterance in the training set with the following objective:

$$\arg \max_{\mathbf{p}_{att}} \mathbb{E}_{\mathbf{u} \sim \mathcal{U}} [\mathcal{L}(B_t, f(\mathbf{p}_{att}; \mathbf{e}_u))],$$

where \mathcal{U} are user utterances and \mathcal{L} is the loss function of the DST task. By maximizing the loss for the ground truth belief states we aim to find prompts that force the model to make the most wrong predictions by pushing far apart from the ground truth, like guessing “expensive” instead of “cheap” for u_t (“I am looking for cheap food.”).

In addition, we explore an alternative tuning objective – minimizing the loss. We replace all the non-empty values in B_t to empty (e.g., (restaurant-price range, expensive) changes to (restaurant-price range, none)) and then minimize the loss:

$$\arg \min_{\mathbf{p}_{att}} \mathbb{E}_{\mathbf{u} \sim \mathcal{U}} [\mathcal{L}(B'_t, f(\mathbf{p}_{att}; \mathbf{e}_u))],$$

where B'_t is the set of target belief states. Different from our previous tuning objective, here we aim to find prompts that force the model to fail to extract the correct value for the slot from user utterances. For example, the DST will fail to extract “cheap” for slot *price range* in u_t (“I am looking for cheap food.”) and thus the predicted belief states will become (restaurant-price range, none).

3.3 Adversarial Example Construction

Next, we focus on the second stage of Prompt-Attack: how can we use these prompts to create adversarial examples that can attack DSTs successfully while maintaining good fluency? After obtaining the adversarial prompts, we use them to generate adversarial examples via mask-and-filling (Li et al., 2021a; Yang et al., 2022; Lei et al., 2022) by pre-trained masked language models. Specifically, we tokenize user utterance u_t to a list of tokens, $u_t = [w_u^1, w_u^2, \dots, w_u^n]$. Then we randomly mask tokens that are not values in B_t , slot-related words, or stopwords with a special token [MASK] and denote the masked u_t as $u_t^m = [w_u^1, \text{[MASK]}, \dots, w_u^n]$. Shown in Figure 2c, we concatenate the adversarial prompts and the masked utterance \mathbf{u}_t^m and use a masked language model \mathcal{M} to predict masked text pieces and generate the perturbations based on surrounded context. As shown in Table 1, for discrete prompt \mathbf{p}_{att}^d , the input for \mathcal{M} would be

Method	$\mathbf{p}_{att} + \mathbf{u}_t^m$ (or \mathbf{e}_u^m)
PromptAttack _d	belief states: $[\mathbf{s}] = [\mathbf{v}]; \mathbf{t}_u^1 [\text{MASK}] \mathbf{t}_u^m$
PromptAttack _c	$\mathbf{p}_1 \mathbf{p}_2 \dots \mathbf{p}_m \oplus \mathbf{e}_u^1 [\text{MASK}] \mathbf{e}_u^m$

Table 1: Adversarial example generation for discrete prompts and continuous prompts.

the concatenation of \mathbf{p}_{att}^d and \mathbf{u}_t^m while for continuous prompt \mathbf{p}_{att}^c , the input would be the concatenation of \mathbf{p}_{att}^c and embedding of masked user utterance $\mathbf{e}_u^1 [\text{MASK}] \mathbf{e}_u^m$. Hence, with \mathbf{p}_{att} and the capability of MLM, the model \mathcal{M} will fill in the blanks with context-consistent tokens which can keep the sentence fluency while maximizing the risk of the DST making wrong predictions, denoted as $P([\text{MASK}] = w | \mathbf{p}_{att}; \mathbf{u}_t^m)$, where w is the generated perturbation. After filling $[\text{MASK}]$ with w and removing \mathbf{p}_{att} , the filled user utterances are used as adversarial examples to attack victim models.

4 Experimental Setup

Our experiments are designed to test the effectiveness of the proposed prompt-based approach to attack DST models. We structure the experiments around four research questions: **RQ1**: Are adversarial examples learned by PromptAttack effective and transferable? And how do these examples compare against baseline (non-prompt) approaches? **RQ2**: Are the generated adversarial examples of good quality? That is, are they fluent with a low perturbation ratio? **RQ3**: What impact do the design choices of PromptAttack have, i.e., the ratio of perturbed tokens and prompt length? **RQ4**: And finally, can the generated adversarial examples be used to improve the performance of current DST models to improve their robustness?

4.1 Dataset

We evaluate our methods on the widely used and challenging multi-domain dialogue dataset, MultiWOZ 2.1 (Eric et al., 2020),² which contains over 10,000 dialogues spanning seven domains. Following existing work (Li et al., 2021b; Lee et al., 2021; Yang et al., 2022), we keep five domains (train, taxi, restaurant, hotel, and attraction) with 30 domain-slot pairs and follow the standard train/validation/test split.

4.2 Evaluation Metrics

We evaluate the proposed methods with a standard set of metrics (Jin et al., 2020; Li et al., 2020,

²github.com/budzianowski/multiwoz, MIT License.

2021a; Simoncini and Spanakis, 2021): **Joint goal accuracy (JGA)**: the average accuracy of predicting all (domain-slot, value) pairs in a turn correctly. **Attack success rate (ASR)**: the proportion of generated adversarial examples that successfully mislead model predictions. **Perturbation ratio (PER)**: the percentage of perturbed tokens in the sentence. Each replace action accounts for one token perturbed. A lower perturbation ratio indicates more semantic consistency (Li et al., 2020). **Perplexity (PPL)**: a metric to evaluate the fluency of sentences. We calculate the perplexity of adversarial examples through GPT-2 (Radford et al., 2019). PPL is calculated across all the adversarial examples. A lower PPL score indicates higher fluency and naturalness of the adversarial examples.

4.3 Baseline Methods

We compare our methods with strong baselines capable of attacking a DST. **TP** and **SD** are two methods maintaining the dialogue act labels unchanged and implemented by the LAUG toolkit (Liu et al., 2021a). For a fair comparison, we do not apply slot value replacement which would modify the slot values in the original utterances. **TP** (Text Paraphrasing) uses SC-GPT (Peng et al., 2020) to generate a new utterance conditioned on the original dialogue acts as data augmentation. **SD** (Speech Disfluency) mimics the disfluency in spoken language by filling pauses (“um”), repeating the previous word, restarting by prepending a prefix “I just” before the original user utterance, and repairing by inserting “sorry, I mean” between a random slot value and the original slot value (Liu et al., 2021a). **SC-EDA** (Liu et al., 2021a) injects word-level perturbations by synonym replacement, random insertions, swaps, and deletions without changing the true belief states. **BERT-M** is introduced in this paper as another baseline method. First, we randomly mask tokens that are not slot-value related and not stopwords. Then, we use BERT (Devlin et al., 2019) to generate perturbations based on the top- K predictions via mask-and-filling, where in our experiments $K = 20$. We sorted the top 20 tokens based on the possibility scores and pick the one with the lowest possibility to fill the masked position. The filled user utterance is regarded as an adversarial example.

4.4 Victim Models

We choose the **TripPy** DST (Heck et al., 2020) as our base model to train our adversarial prompts

Method	TripPy	CoCo	SimpleTOD	TRADE
	JGA↓ / Δ / ASR↑	JGA↓ / Δ / ASR↑	JGA↓ / Δ / ASR↑	JGA↓ / Δ / ASR↑
Original	61.3 / - / -	62.6 / - / -	56.0 / - / -	49.4 / - / -
SC-EDA	60.5 / -0.8 / 1.9	61.9 / -0.7 / 1.6	53.6 / -2.4 / 9.5	48.8 / -0.6 / 4.9
TP	60.3 / -1.0 / 5.6	61.5 / -1.1 / 4.7	52.6 / -3.4 / 19.3	48.8 / -0.6 / 14.1
SD*	56.5 / -4.8 / 9.3	56.1 / -6.5 / 11.4	38.8 / -17.2 / 36.6	31.7 / -17.7 / 39.9
BERT-M	58.9 / -2.4 / 5.0	60.1 / -2.5 / 4.8	49.6 / -6.4 / 16.4	45.9 / -3.5 / 11.5
PromptAttack _d	53.6 / -7.7 / 16.0	53.7 / -8.9 / 16.9	38.9 / -17.1 / 37.9	35.8 / -13.6 / 34.0
PromptAttack _{cx}	53.3 / -8.0 / 16.3	54.1 / -8.5 / 16.3	25.0 / -31.0 / 60.0	35.7 / -13.7 / 34.1
PromptAttack _{cn}	52.0 / -9.3 / 18.2	52.8 / -9.8 / 18.4	37.4 / -18.6 / 40.6	35.8 / -13.6 / 33.3

Table 2: Attack effectiveness results on MultiWOZ 2.1. **JGA (%)**: joint goal accuracy; **Δ (%)**: the absolute difference between original JGA and JGA after attacking; **ASR (%)**: attack success rate. ↓ (↑): denotes whether the lower (or higher) the better from an attack perspective. *: denotes the method that introduces new slot values. We highlight the **best** and the second best results.

since classification-based models have better performance and are more robust than generation-based models (Liu et al., 2021a). Demonstrating the susceptibility of TripPy to our adversarial examples can reveal the limitations of current DSTs, but we further explore the *transferability* of the prompt-based attacks.

Transferability reflects the generalization of the attack methods, meaning that adversarial examples generated for one model can also effectively attack other models (Zhang et al., 2020). Hence, we also evaluate the prompt-based approach learned over TripPy by targeting our adversarial examples on other popular DSTs: **TRADE** (Wu et al., 2019), **SimpleTOD** (Hosseini-Asl et al., 2020), and **CoCo** (Li et al., 2021b), one of the state-of-the-art models.³ Additional information about the implementations can be found in Appendix A.

5 Experimental Results

Given this setup, we now investigate the four experimental research questions in turn.

5.1 Attack Effectiveness (RQ1)

First, are the adversarial examples learned by PromptAttack effective? Table 2 summarizes the results for three versions of PromptAttack versus the baselines for the four different DSTs (TripPy, CoCo, SimpleTOD, and TRADE). We consider the discrete version of PromptAttack (denoted as **PromptAttack_d**) and two continuous versions: one is optimized by maximizing the training loss (denoted as **PromptAttack_{cx}**), while the other one is optimized by minimizing the loss (denoted as **PromptAttack_{cn}**).

³These models are fine-tuned on MultiWOZ 2.1 using code from CoCo (<https://github.com/salesforce/coco-dst>) and follow the same post-processing strategy as CoCo. BSD 3-Clause License.

Attack Performance. First, let’s focus on the TripPy column. All versions of PromptAttack are learned over TripPy and then applied here so we can assess the susceptibility of a popular DST to adversarial examples. The four baselines lead to some degradation in terms of accuracy (JGA), with SD performing the best with a JGA of 56.5 (a 4.8 drop from the original DST).⁴ The three prompt-based learning approaches result in strong degradation in terms of accuracy, ranging from 7.7 to 9.3 drops relative to the original. We observe that our PromptAttack models significantly outperform SC-EDA, TP, and BERT-M, the methods without introducing new slot values in the adversarial examples, in terms of JGA and ASR. Compared with the best baseline method among these three, BERT-M, PromptAttack_{cn} decreases the JGA by 6.9 and increases ASR by 13.2, respectively. In addition, for the method introducing new slot values, SD, PromptAttack_{cn} outperforms it by 4.5 and 8.9. Hence, these observations reveal the attack effectiveness of our proposed PromptAttack methods over these baselines no matter whether the methods introduce new slot values or not.

Transferability. To test the transferability of the generated adversarial examples, we take the examples trained over TripPy and then use them to attack other victim models CoCo, SimpleTOD, and TRADE. For CoCo and SimpleTOD, we see that PromptAttack outperforms these four baselines. Our best method PromptAttack_c achieves

⁴We attribute this good attack performance since although this method maintains ground truth slot-value labels unchanged, it prepends new slot values before the original slot values in the user utterance. This operation is effective because it can easily confuse the model to decide which slot values are the truth slot values. In contrast, our prompt-based approaches are designed to make very few changes and to avoid introducing new slot values.

52.8 and 25.0 JGA when attacking CoCo and SimpleTOD, showing better transferability than PromptAttack_d. For TRADE, PromptAttack_c shows better attack performance than baselines without introducing new slot values significantly. Specifically, PromptAttack_{cx} shows a decrease of 10.2 and an increase of 20.0 in terms of JGA and ASR, respectively. In general, our PromptAttack methods show good transferability: the adversarial examples generated for one victim model can also be used to attack another model effectively.

5.2 Adversarial Example Quality (RQ2)

Next, we examine whether the generated adversarial examples are of good quality. First, are they fluent with a low perturbation ratio? We automatically measure the perturbation ratio (PER) between the original input and adversarial examples, and the fluency by computing the perplexity (PPL). The lower perturbation ratio represents fewer perturbed tokens in original utterances and lower perplexity indicates better fluency. From Table 3 we observe that the PromptAttack methods achieve low perplexity and show good fluency with quite a low perturbation ratio. Specifically, our method PromptAttack_{cn} (7.7%) achieves 169.0 PPL, showing better fluency than PromptAttack_{cn} (28.1%) and baselines. Although SC-EDA has a lower perturbation ratio than our PromptAttack_{cn} (28.1%), it shows less attack effectiveness (Section 5.1) and worse fluency. Thus, there are trade-offs between perturbation ratio and attack effectiveness.

Second, do the adversarial examples preserve the semantics of the un-perturbed original sentences? That is, does an utterance asking for a cheap restau-

Method	PER↓	PPL↓	Semantic↑	Grammar↑
Original	-	173.7	-	4.8
SC-EDA	13.1	773.8	2.5	2.7
TP	74.4 [†]	352.4	2.6	4.8
SD*	30.4 [†]	270.4	4.3	4.1
BERT-M	28.1	221.3	2.8	4.3
Adv (7.7%)	7.7	169.0	4.3	<u>4.4</u>
Adv (28.1%)	28.1	<u>177.6</u>	<u>3.3</u>	3.8

Table 3: Automatic evaluation and human evaluation results. **PER**: perturbation ratio; **PPL**: perplexity of generated adversarial examples representing fluency. ↓ (↑) denotes whether the lower (or higher) is the better. [†]: results are from original papers. * denotes the method that introduces new slot values. Adv (*): adversarial examples from PromptAttack_{cn} with different perturbation ratios which lead the victim model’s accuracy to 0. We highlight the **best** and the second best results.

rant lead to an adversarial example that also asks for a cheap restaurant though tricking the DST to output expensive? To answer this question, we conduct a human evaluation on semantics preservation and grammatical correctness. We first shuffled 150 examples: 50 original un-perturbed sentences, 50 adversarial examples with a 7.7% perturbation ratio, and 50 with a 28.1% perturbation ratio (following the analysis in Section 5.3.1). For the adversarial examples, each attacks the victim model successfully leading to an accuracy of 0. Following (Jin et al., 2020; Li et al., 2020), we ask three human judges to rate how well a randomly chosen sentence preserves the semantics of the original sentence (*semantic*), how grammatically correct the sentence is (*grammar*), on a scale from 1 to 5. We report the average score across the three judges in Table 3.

As we can see, the semantic score and grammar score of the adversarial examples are close to the original ones. We find that when the perturbation is reasonable (around 8%), the semantics of the original sentence are preserved quite well (scoring 4.3 for adversarial examples). Further, the grammatical quality of the sentence is also maintained well (4.8 versus 4.4). Even as the perturbation ratio increases to approximately 28%, our approach continues to uphold good semantic preservation (3.3) while retaining satisfactory grammar quality (3.8). Overall, our method consistently generates high-quality adversarial examples by effectively preserving semantics, maintaining grammatical quality and fluency, and keeping a low perturbation ratio.

5.3 Impact of PromptAttack Design (RQ3)

We now explore the impact of different settings on our proposed methods.

5.3.1 Ratio of Perturbed Tokens

First, our prompt-based approach can control how many tokens we want to change in the user utterances, which gives it flexibility. Since the perturbation ratio represents the semantic consistency between the original examples and adversarial examples and there are trade-offs between the attack effectiveness and perturbation ratio, it is important to investigate the influence of the ratio of perturbed tokens on attacking ability.

We take $\max(1, \text{perturbation_ratio} * l_t)$ as the number of perturbed tokens, where l_t denotes the length of pre-processed utterances. We set the perturbation ratio of tokens that we could perturb to

		7.7%	10.2%	15.2%	22.6%	28.1%
		(1.0)	(1.5)	(2.3)	(3.5)	(4.4)
P_{cx}	JGA↓	59.0	58.1	56.6	55.1	53.3
	ASR↑	4.6	6.3	9.5	12.8	16.3
	PPL↓	159.4	155.9	157.5	167.0	175.5
P_{cn}	JGA↓	58.9	58.1	56.4	54.0	52.0
	ASR↑	4.9	6.2	9.7	14.4	18.2
	PPL↓	169.0	173.7	177.1	172.2	177.6

Table 4: Results of PromptAttack_{cx} (P_{cx}) and PromptAttack_{cn} (P_{cn}) with different perturbation ratio. (*) denotes the average perturbed token numbers.

10%, 30%, 50%, 80%, and 100%, that is 7.7%, 10.2%, 15.2%, 22.6%, and 28.1% of the average length of all input examples. More data analysis can be found in Appendix B.

Table 4 shows the evaluation of attack performance and fluency of generated adversarial examples from PromptAttack_{cx} and PromptAttack_{cn}. We observe that for these two methods, the more tokens we perturb, the lower JGA and higher ASR we get, showing better attack ability, which is consistent with our intuition. Thus, as the ratio of perturbed tokens increases, our proposed method PromptAttack achieves better attack performance while maintaining good fluency.

5.3.2 Prompt Length

Next, we explore the effect of different continuous prompt lengths. Shorter prompts have fewer tunable parameters, which means under the same training setting, it would be faster to optimize and find the most effective adversarial prompts. We train continuous prompts with different length: 5 tokens, 10 tokens, and 15 tokens using PromptAttack_{cx}. Table 5 shows that under different prompt lengths, with the increase of perturbation ratio, the model achieves better attack performance. Under the same perturbation ratios, the model with 5-token prompt achieves modest lower JGA and higher ASR. For example, when the perturbation ratio is 28.1%, PromptAttack_{cx} with 5-token prompt gains lower JGA than PromptAttack_{cx} with 10-token prompt

	P_5		P_{10}		P_{15}	
	JGA↓	ASR↑	JGA↓	ASR↑	JGA↓	ASR↑
7.7%	59.0	4.6	59.2	4.3	59.3	4.6
10.2%	58.1	6.3	58.5	5.9	58.5	6.0
15.2%	56.6	9.5	57.0	8.8	57.2	8.8
22.6%	55.1	12.8	55.3	12.6	55.7	12.3
28.1%	53.3	16.3	53.5	15.9	54.1	15.1

Table 5: Results of PromptAttack_{cx} with different prompts length and perturbation ratios. P_* denotes the prompt length.

	JGA _d	JGA _o	ASR _d	ASR _o
Original	67.3	61.3	-	-
SC-EDA	66.5	60.5	1.8	1.9
TP	65.9	60.3	5.5	5.6
SD*	61.4	56.5	10.1	9.3
BERT-M	64.5	58.9	5.0	5.0
PromptAttack _d	60.0	55.8	12.6	11.3
PromptAttack _{cx}	58.3	53.3	16.3	16.3
PromptAttack _{cn}	56.8	52.0	18.5	18.2

Table 6: Defense results. d: defended DST model; o: original DST model.

and PromptAttack_{cx} with 15-token prompt by 0.2 and 0.8, respectively, and higher ASR by 0.4 and 1.2, indicating slightly better attack performance.

5.4 Defense against Attack (RQ4)

Finally, we turn to the challenge of defending a DST in the presence of such adversarial examples. We aim to answer two questions: i) can our generated adversarial examples be used to improve the performance of current DST models? and ii) can our attack method bypass such a defense method?

One of the most effective approaches to increase the robustness of a model is adversarial training, which injects adversarial examples into the training data to increase model robustness intrinsically (Bai et al., 2021). Specifically, we first apply our attack methods on the original training dataset to generate adversarial examples. Then we re-train the TripPy model on the training set augmented by the adversarial training examples and evaluate the performance on original test set. As shown in Table 6, the new defended DST model improves JGA on the original test set from 61.3 to 67.3 by 6.0, which outperforms results reported by the state-of-the-art DST model CoCo (62.6) by 4.7. This encouraging result shows that adversarial examples from our attack method can be a good source for data augmentation.

To evaluate the robustness of such an augmented DST model against our proposed attack methods, we next test how well our adversarial examples perform. From Table 6 we observe that the attack methods still show strong attack ability on the new DST model. Thus, there is an opportunity to explore stronger defense methods to strengthen DSTs against such prompt-based attacks.

6 Conclusion

In this paper, we present a prompt-based learning approach that can generate effective adversarial

examples for probing DST models. Through experiments over four state-of-the-art DSTs, our framework achieves the greatest reduction in accuracy with the best attack success rate. Moreover, the generated adversarial examples maintain good fluency and low perturbation ratio, evidence that they are close to legitimate non-adversarial user inputs. We also show our generated adversarial examples can bolster a DST by augmenting the original training data with adversarial examples. We find that both discrete and continuous adversarial prompts are capable of generating effective adversarial examples. Discrete prompts are more interpretable while continuous prompting allows us to search for optimal adversarial prompts more efficiently, and generates more effective adversarial examples.

Limitations

The natural idea to improve robustness is to add adversarial examples to the training set and retrain the model. However, generating adversarial examples for a large training set can be very time-consuming. Thus, it would be interesting to explore more efficient methods that implicitly involved adversarial examples in the training process, e.g., (Yang et al., 2022).

Ethics Statement

The proposed methods could also be applied to natural language generation tasks, like dialogue response generation. The misuse of such methods may generate biased or offensive responses.

Acknowledgements

We appreciate the authors of the CoCo paper for making their code accessible to the public and for taking the time to address our inquiries.

References

Tao Bai, Jinqi Luo, Jun Zhao, Bihan Wen, and Qian Wang. 2021. [Recent advances in adversarial training for adversarial robustness](#). In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 4312–4321. International Joint Conferences on Artificial Intelligence Organization. Survey Track.

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

and Dialogue, pages 239–251, Singapore and Online. Association for Computational Linguistics.

- Max Bartolo, Tristan Thrush, Robin Jia, Sebastian Riedel, Pontus Stenetorp, and Douwe Kiela. 2021. [Improving question answering model robustness with synthetic adversarial data generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8830–8848, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yonatan Belinkov and Yonatan Bisk. 2018. Synthetic and natural noise both break neural machine translation. In *International Conference on Learning Representations*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- 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 dialogue state tracking](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 879–885, Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- 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, Marseille, France. European Language Resources Association.
- Yue Feng, Yang Wang, and Hang Li. 2021. [A sequence-to-sequence approach to dialogue state tracking](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th*

- International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1714–1725, Online. Association for Computational Linguistics.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3816–3830, Online. Association for Computational Linguistics.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. *ICLR 2015*.
- Yun He, Steven Zheng, Yi Tay, Jai Gupta, Yu Du, Vamsi Aribandi, Zhe Zhao, YaGuang Li, Zhao Chen, Donald Metzler, et al. 2022. Hyperprompt: Prompt-based task-conditioning of transformers. In *International Conference on Machine Learning*, pages 8678–8690. PMLR.
- Michael Heck, Carel van Niekerk, Nurul Lubis, Christian Geisshauser, 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, 1st virtual meeting. Association for Computational Linguistics.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. *Advances in Neural Information Processing Systems*, 33:20179–20191.
- Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2021–2031.
- Di Jin, Zhijing Jin, Joey Zhou, and Peter Szolovits. 2020. [Is bert really robust? a strong baseline for natural language attack on text classification and entailment](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:8018–8025.
- Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2021. [Dialogue state tracking with a language model using schema-driven prompting](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4937–4949, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yibin Lei, Yu Cao, Dianqi Li, Tianyi Zhou, Meng Fang, and Mykola Pechenizkiy. 2022. Phrase-level textual adversarial attack with label preservation. *NAACL-HLT 2022 Findings*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Dianqi Li, Yizhe Zhang, Hao Peng, Liqun Chen, Chris Brockett, Ming-Ting Sun, and Bill Dolan. 2021a. [Contextualized perturbation for textual adversarial attack](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5053–5069, Online. Association for Computational Linguistics.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. [BERT-ATTACK: Adversarial attack against BERT using BERT](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202, Online. Association for Computational Linguistics.
- Shiyang Li, Semih Yavuz, Kazuma Hashimoto, Jia Li, Tong Niu, Nazneen Rajani, Xifeng Yan, Yingbo Zhou, and Caiming Xiong. 2021b. [Coco: Controllable counterfactuals for evaluating dialogue state trackers](#). In *International Conference on Learning Representations*.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.
- Jiexi Liu, Ryuichi Takanobu, Jiaxin Wen, Dazhen Wan, Hongguang Li, Weiran Nie, Cheng Li, Wei Peng, and Minlie Huang. 2021a. [Robustness testing of language understanding in task-oriented dialog](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2467–2480, Online. Association for Computational Linguistics.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021b. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.
- Jiaying Lu, Jiaming Shen, Bo Xiong, Wenjing Ma, Stefan Staab, and Carl Yang. 2023. Hiprompt: Few-shot biomedical knowledge fusion via hierarchy-oriented prompting. *arXiv preprint arXiv:2304.05973*.
- Andrea Madotto, Zhaojiang Lin, Genta Indra Winata, and Pascale Fung. 2021. Few-shot bot: Prompt-

- based learning for dialogue systems. *arXiv preprint arXiv:2110.08118*.
- Shikib Mehri, Mihail Eric, and Dilek Hakkani-Tur. 2020. Dialoglue: A natural language understanding benchmark for task-oriented dialogue. *arXiv preprint arXiv:2009.13570*.
- John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. **TextAttack: A framework for adversarial attacks, data augmentation, and adversarial training in NLP**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 119–126, Online. Association for Computational Linguistics.
- Baolin Peng, Chenguang Zhu, Chunyuan Li, Xiujuan Li, Jinchao Li, Michael Zeng, and Jianfeng Gao. 2020. **Few-shot natural language generation for task-oriented dialog**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 172–182, Online. Association for Computational Linguistics.
- Danish Pruthi, Bhuwan Dhingra, and Zachary C. Lipton. 2019. **Combating adversarial misspellings with robust word recognition**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5582–5591, Florence, Italy. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. **Generating natural language adversarial examples through probability weighted word saliency**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1085–1097, Florence, Italy. Association for Computational Linguistics.
- Walter Simoncini and Gerasimos Spanakis. 2021. **SeqAttack: On adversarial attacks for named entity recognition**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 308–318, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yibo Wang, Congying Xia, Guan Wang, and Philip S. Yu. 2022. **Continuous prompt tuning based textual entailment model for e-commerce entity typing**. In *2022 IEEE International Conference on Big Data (Big Data)*, pages 1383–1388.
- 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 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 808–819.
- Yuting Yang, Pei Huang, Juan Cao, Jintao Li, Yun Lin, Jin Song Dong, Feifei Ma, and Jian Zhang. 2022. **A prompting-based approach for adversarial example generation and robustness enhancement**. *arXiv preprint arXiv:2203.10714*.
- Yuting Yang, Wenqiang Lei, Pei Huang, Juan Cao, Jintao Li, and Tat-Seng Chua. 2023. **A dual prompt learning framework for few-shot dialogue state tracking**.
- Tao Yu, Rui Zhang, Alex Polozov, Christopher Meek, and Ahmed Hassan Awadallah. 2021. **{SC}ore: Pre-training for context representation in conversational semantic parsing**. In *International Conference on Learning Representations*.
- Wei Emma Zhang, Quan Z Sheng, Ahoud Alhazmi, and Chenliang Li. 2020. Adversarial attacks on deep-learning models in natural language processing: A survey. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(3):1–41.
- Jeffrey Zhao, Mahdis Mahdieh, Ye Zhang, Yuan Cao, and Yonghui Wu. 2021. **Effective sequence-to-sequence dialogue state tracking**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7486–7493, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021. **Factual probing is [MASK]: Learning vs. learning to recall**. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5017–5033, Online. Association for Computational Linguistics.
- Qi Zhu, Bing Li, Fei Mi, Xiaoyan Zhu, and Minlie Huang. 2022. **Continual prompt tuning for dialog state tracking**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1124–1137, Dublin, Ireland. Association for Computational Linguistics.

A Implementation Details

A.1 Constructing Adversarial Examples

We use the TripPy DST model (Heck et al., 2020) as the victim model, which uses the 12-layer pre-trained BERT-base-uncased model (Devlin et al., 2019) as the context encoder and has 768 hidden units, 12 self-attention heads, and 110M parameters. We train our adversarial prompts for 10 epochs with an initial learning rate of 1×10^{-4} . The LR decay linearly with a warmup proportion of 0.1. We use Adam optimizer for optimization and set the maximum input sequence length of user utterance l_u to 180 and the number of prompt tokens l_p to $\{5, 10, 15\}$. The total length of the input is

$l_u + l_p$. The training batch size is 64 and the evaluation batch size is 1. We evaluate the checkpoint of the prompt for each epoch and choose the one that leads to the lowest JGA on the validation set as our final adversarial prompt. The MLM model used to generate the adversarial examples via mask-and-filling is also BERT-base-uncased.

A.2 Training Defense Models

We train the TripPy defense DST model on the training dataset augmented with adversarial examples following the training setting in (Heck et al., 2020). The model uses the pre-trained BERT-base-uncased transformer as the context encoder front-end, which has 12 hidden layers with 768 hidden units and 12 self-attention heads each (Heck et al., 2020). The initial learning rate is set to 1×10^{-4} with a warmup proportion of 0.1 and let the LR decay linearly after the warmup phase. We use Adam optimizer for optimization and dropout on the BERT output with a rate of 30%. The training batch size is 48 and the model is trained for 10 epochs with early stopping employed based on the JGA on the validation set. The experiments are run on 2 NVIDIA TITAN Xp GPUs.

B Data Analysis

Figure 3 shows the distribution of length of original user utterance l_o , the length of utterances after removing stop words, slot and value related tokens l_t , and the difference between them, that is $\Delta = l_o - l_t$. We can see, 95.5% of user utterances have fewer than 10 tokens that could be perturbed and 59.3% of them could perturb less than 4 tokens.

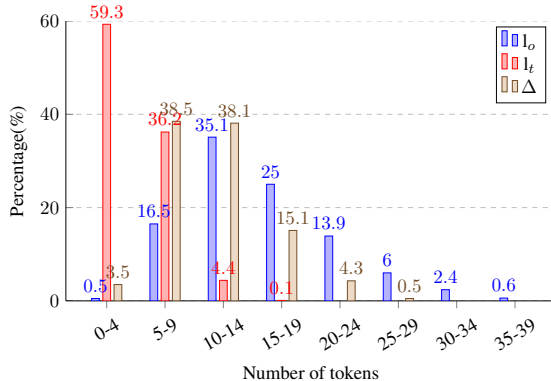


Figure 3: Data analysis. l_o : original length of user utterance; l_t : after data pre-processing, the number of available tokens that can be perturbed; Δ : $l_o - l_t$.

C Success Generation Rates of Adversarial Examples

To assess whether there are any masked tokens kept unchanged after the adversarial example construction process, we evaluate the success generation rates of adversarial examples generated using our methods. we observe that PromptAttack_d achieves a success rate of 0.94, PromptAttack_{cx} achieves 0.95, and PromptAttack_{cn} achieves 0.96. These findings indicate that our methods demonstrate a high rate of successful generation. And it is important to note that unchanged utterances will not attack the models successfully.

D Data Augmentation

In Section 5.4, TripPy trained on original training data augmented with our generated adversarial examples improves TripPy by 6.0% and outperforms CoCo by 4.7% when evaluated on the original test set following the same post-processing strategy as CoCo (Li et al., 2021b). When using TripPy’s default cleaning, the comparison results with previous methods are shown in Table 7.

Model	JGA (%)
TRADE (Wu et al., 2019)	46.00 [†]
TripPy (Heck et al., 2020)	55.29 [†]
SimpleTOD (Hosseini-Asl et al., 2020)	55.76 [†]
ConvBERT-DG + Multi (Mehri et al., 2020)	58.70 [†]
TripPy + SCoRE (Yu et al., 2021)	60.48 [†]
TripPy + CoCo (Li et al., 2021b)	60.53 [†]
Ours	60.56
TripPy + SaCLog (Dai et al., 2021)	60.61 [†]

Table 7: DST results on MultiWOZ 2.1. [†] denotes that results are from original papers.

E Case Study

In Table 8, we showcase some examples of the adversarial examples generated from baseline methods and our methods.

F Error Analysis

Although we keep the ground-truth values and slot-related words unchanged in the adversarial examples, perturbations on surrounding words may lead to low human evaluation results. Here we show some adversarial examples that attack the victim model successfully in Table 9. For case 1, changing “night” to “hour” affects the meaning of the sentence but is consistent with ground truth “1”. However, the victim model predicts “two” for slot

Methods	User Utterance	Predictions	Accuracy
Ori	Can you give me the address of the adc theatre please?	attraction-name-adc theatre	1.0
LAUG-SD	Oh , can you give me the address of the uh, galleria nope , adc theatre please?	attraction-name-none; restaurant-name-galleria	0.0
LAUG-TP	Adc theatre sounds good. What is the address?	attraction-name-none	0.0
BERT-M	Can you email me the address of the adc theatre please?	attraction-name-adc theatre	1.0
PromptAttack _d	Can you give me the name of the adc theatre please?	attraction-name-adc theatre	1.0
PromptAttack _c	Can you give me the support of the adc theatre please?	attraction-name-none	0.0

Table 8: Cases of adversarial examples generated from different attacking methods. **Bold** indicates the difference from original user utterance.

“book_stay”. For cases 2 and 3, we could see by modifying words “well” and “restaurant”, which are not related to ground-truth slot-values, the victim model can not identify “six” and “entertainment” for slot “book_people” and “attraction-type”. For cases 4, 6, and 9, the adversarial examples modified words “star”, “leaves”, and “price range” that are related to slot names “stars”, “departure”, and “pricerange” and cause victim models can not identify the corresponding values. In future work, it’s important to consider more comprehensively to prevent modifying slot-related words. For case 5, the new words “office” confuse the model and make it ignore the values “center” and predict “none” for slot “area”. For cases 7, 9, and 10, “fishing”, “french english”, and “street” in adversarial examples lead victim models prediction new value “fishing boat”, “french”, and “cambridge street” for slots “attraction-type”, “restaurant-food”, and “train-departure”. For case 8, the remove of “arriving” leads victim model predict “none” for “train-arriveBy”.

No.	Original Example	Adversarial Example	Ground Truth	Predictions
1	How about 1 night ? Would that work?	How about 1 hour ? Would that work?	hotel-book_stay-1	hotel-book_stay-two
2	The restaurant is for six as well .	The restaurant is for six as t .	restaurant-book_people-6	restaurant-book_people-none
3	I'm also looking for some entertainment close to the restaurant . Any suggestions?	I'm also looking for some entertainment close to the internet . Any suggestions?	attraction-type-entertainment	attraction-type-none
4	Can you find me a three star place to stay?	Can you find me a three some place to stay?	hotel-stars-3	hotel-stars-none
5	What's the phone number for the one in the center?	What's the phone number for the office in the center?	hotel-area-center	hotel-area-none
6	I'm also looking for a train that leaves leicester.	I'm also looking for a train that is leicester.	train-departure-leicester	train-departure-leicester, train-destination-leicester
7	Any interesting boats on the east side of town?	Any fishing boats on the east side of town?	attraction-type-boats	attraction-type-fishing-boats
8	I also need a train to Cambridge arriving at 10:15 on Thursday.	I will need a train to Cambridge university at 10:15 on Thursday.	train-arriveBy-10 : 15, train-day-thursday, train-destination-cambridge	train-leaveAt-10 : 15, train-arriveBy-none, train-day-thursday, train-destination-cambridge
9	I'm looking for a restaurant in the north part and in cheap price range .	I'm looking for a restaurant in the north part and in cheap french english .	restaurant-area-north, restaurant-pricerange-cheap	restaurant-area-north, restaurant-food-french, restaurant-pricerange-cheap
10	Can you help me find a train departing from cambridge going to kings lynn?	Can you help me find a train departing from cambridge street to kings lynn?	train-departure-cambridge, train-destination-kings lynn	train-departure-cambridge street, train-destination-kings lynn

Table 9: Examples for error analysis.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Limitations
- A2. Did you discuss any potential risks of your work?
Ethics Statement
- A3. Do the abstract and introduction summarize the paper's main claims?
Abstract, Section 1 Introduction
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

Section 4 Experimental Setup

- B1. Did you cite the creators of artifacts you used?
Section 4 Experimental Setup
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
Section 4 Experimental Setup
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
Section 4 Experimental Setup
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
Not applicable. Left blank.
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
Section 4.1 Dataset
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
Section 4.1 Dataset, Appendix B

C Did you run computational experiments?

Section 4 Experimental Setup

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
Appendix A

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

Appendix A

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

Appendix A

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Appendix A

D Did you use human annotators (e.g., crowdworkers) or research with human participants?

Section 5.2 Adversarial Example Quality

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

We report the text of instructions about two metrics (Semantic and Grammar) which need to be annotated in Section 5.2 Adversarial Example Quality. There are no examples of risks in this study.

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

The annotators were graduate students in the computer science area and they were not compensated.

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

The annotators knew the data would be used for this paper.

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

Not applicable. Left blank.

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

Not applicable. Left blank.