

---

# Say What I Want: Towards the Dark Side of Neural Dialogue Models

---

**Haochen Liu**

Michigan State University  
liuhaoc1@msu.edu

**Tyler Derr**

Michigan State University  
derrtyl@msu.edu

**Zitao Liu**

TAL AI Lab  
liuzitao@100tal.com

**Jiliang Tang**

Michigan State University  
tangjili@msu.edu

## Abstract

Neural dialogue models have been widely adopted in various chatbot applications because of their good performance in simulating and generalizing human conversations. However, there exists a dark side of these models – due to the vulnerability of neural networks, a neural dialogue model can be manipulated by users to say what they want, which brings in concerns about the security of practical chatbot services. In this work, we investigate whether we can craft inputs that lead a well-trained black-box neural dialogue model to generate targeted outputs. We formulate this as a reinforcement learning (RL) problem and train a **Reverse Dialogue Generator** which efficiently finds such inputs for targeted outputs. Experiments conducted on a representative neural dialogue model show that our proposed model is able to discover such desired inputs in a considerable portion of cases. Overall, our work reveals this weakness of neural dialogue models and may prompt further researches of developing corresponding solutions to avoid it.

## 1 Introduction

Dialogue system, also known as conversational AI, which aims to conduct human-like conversations with users, is receiving increasing attention from both the industry and the academic research community. In the past, such systems either rely on intricate hand-crafted rules [1, 2], or depend on a complicated processing pipeline including a series of functional modules [3]. Meanwhile, retrieval-based methods [4, 5, 6, 7], which search a suitable response from a repository given the query, are also adopted in many application scenarios. These methods are able to provide natural, human-like responses, but fail to generate novel responses out of the range of the repository [8]. Recently, researchers begin to involve deep learning techniques in building fully data-driven and end-to-end dialogue systems [9], which are referred to as neural dialogue models.

Based on the Seq2Seq framework [10], these neural models [11, 12, 13, 14, 15] have achieved surprising performances and gradually dominate the field of dialogue generation. First, these models are easy to train. Instead of designing complicated rules or modularized pipelines, the models can learn the mapping between queries and responses automatically from massive existing dialogue pairs [12]. Second, given that the neural models are trained on large-scale human conversation data, they show a strong generalization ability that they can handle open-domain conversations rather than restricting the topics in a narrow domain [12, 16]. In addition, neural dialogue models can provide fluent and smooth responses, and show the intelligence of performing simple common sense reasoning [12]. Since neural dialogue models achieve a breakthrough of conducting reasonable and

engaging human-like conversations, they are widely adopted by the industry as a core component of practical chatbot applications, such as Microsoft XiaoIce [16].

While the research community is delighted with the success of neural dialogue models, there is a dark side of these models. Given that the internal mechanisms of neural networks are not explicitly interpretable, neural dialogue models are vulnerable. For example, they may have unpredictable behaviors with regard to some well-crafted inputs [17, 18]. This vulnerability can cause a series of problems, one of which is, whether we can manipulate a dialogue agent to say what we want. In other words, can we find a well-designed input, to induce the agent to provide a desired output? If this is possible, people with ulterior motives may take advantage of this weakness of the chatbots to guide them say something malicious or sensitive, causing adverse social impacts [19, 20, 21].

In this paper, we want to study this dark side by seeking an answer to the question – whether we can design an algorithm that can automatically generate inputs that lead a state-of-the-art black-box neural dialogue model to “say what I want”. However, this presents tremendous challenges. First, unlike similar works such as [22], where the authors try to craft inputs for a neural image captioning model to output targeted sentences, our problem involves discrete inputs (i.e. texts rather than images) and treats the model as a black-box (since the setting is more realistic). Thus, the traditional optimization method that finds the inputs by the guidance of gradient information is completely invalid. Second, when trying to manipulate a dialogue system released by others, it is impractical for us to interact with it for unlimited times. Based on this point, brute-force search methods cannot be adopted, and the number of the interactions with the black-box model that we need to find an input for a targeted output should be restricted to a reasonable level.

To address the above challenges, for a given black-box neural dialogue model, we propose to train a corresponding Reverse Dialogue Generator, which takes a targeted response as input and automatically outputs a query that leads the dialogue model to return that response. The proposed Reverse Dialogue Generator is based on a Seq2Seq model and performs as a reinforcement learning (RL) agent. The black-box dialogue model is regarded as the environment the agent interacts with. It is optimized through policy gradients, with the similarity between the targeted outputs and what the dialogue model outputs with regard to a crafted input as the reward signal. Extensive experiments conducted on a public well-trained neural dialogue model demonstrate the capacity of our model.

## 2 Related Works

Basically, our work is related to the problem of model attacks. Although deep learning models have been used for many tasks that have shown to be useful across a plethora of domains, more recently, researchers have become aware to the fact that although these systems perform extremely well when in a perfect and stable environment (the type they were designed in), but when placed in the real world, they are quite easily susceptible to being attacked. Szegedy et al. [17] first investigates the vulnerability of DNN-based image classifiers by crafting adversarial examples with imperceptible perturbations that lead the classifier to make mistakes. Besides, Sharif et al. [23] focus on attacking face recognition models; Xie et al. [24] try to find the weakness of an object detection system; in [25, 26], researchers study the robustness of generative models, and novel methods to make adversarial examples for deep reinforcement learning based models are introduced in [27, 28].

Adversarial attacks on deep learning models for NLP tasks also attract a lot of interest. Attacking NLP models are more challenging since the inputs are discrete texts instead of continuous values such as image inputs [29]. Text classification problems are studied in [30, 31]; sentiment analysis is involved in [32]; grammar error detection is investigated in [33]. And Belinkov et al. [34] try to fool a machine translation system, while Chan et al. [35] study how to attack a neural reading comprehension model. Besides, more works can be found in the survey [29].

For dialogue generation task, Wieting et al. [36] explore the over-sensitivity and over-stability of neural dialogue models by using some heuristic techniques to modify original inputs and observe the corresponding outputs. They evaluate the robustness of dialogue models by checking whether the outputs change significantly after the modifications on the inputs but don’t consider targeted outputs. The work which is most related to ours is [21], where the authors try to find trigger inputs which can lead a neural dialogue model to output a list of targeted egregious responses. Different from our work, they treat the dialogue model as a white-box and take advantage of the model structure and parameters. By comparison, our black-box setting is more challenging. Furthermore, their algorithm

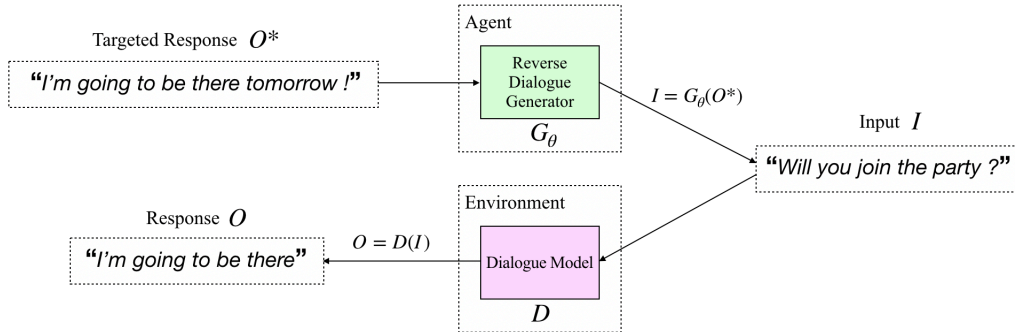


Figure 1: The agent-environment setup of the proposed framework.

indeed fails to lead the dialogue model to output the exact targeted responses. It should be pointed out that similar with [21], our work is not a model attack task in the true sense since we only focus on the requirements of the outputs but don't force the inputs to be close to any original inputs. However, our problems should be solved following the same ideas as adversarial attack problems.

By the way, nevertheless Cheng et al. [37] don't focus on dialogue generation, they also try to manipulate a seq2seq model to generate texts with certain restrictions. However, their non-overlapping and targeted keywords settings are looser than ours where a whole targeted response is required to be output, and this work is under the white-box assumption.

### 3 Reverse Dialogue Generator

In this work, we consider the specific neural dialogue model  $D$  of our interest as a black-box environment. The dialogue model  $D$  is able to take an input sentence  $I$  and output a corresponding response sentence  $O$ . Now given a targeted output  $O^*$ , our goal is to find a well-designed input  $I$ , which leads the dialogue model  $D$  to output a response  $O$  that is exactly the same as  $O^*$  ( $O = O^*$ ), or at least similar with  $O^*$  in semantics ( $O \approx O^*$ ). To achieve it, we build a Reverse Dialogue Generator agent  $G_\theta$ , which takes the targeted output  $O^*$  as input to predict its corresponding input  $I$ . A sketch of this agent-environment setup is shown in Figure 1.

#### 3.1 The Dialogue Model Environment

In this work, we adopt the classical Seq2Seq neural dialogue model [12] as the dialogue model environment. In the model, the encoder and the decoder are implemented by two independent RNN structures or its variants such as LSTM [38] and GRU [39]. The encoder reads the input sentence word by word and encodes it as a low-dimensional latent representation, which is then fed into the decoder to predict the corresponding response sentence word by word. We assume the environment is a black-box; once the model is well-trained, we keep it sealed and the agent has no access to its structures, parameters or gradients.

Thanks to its strong power in the simulation and generalization of human conversations, Seq2Seq-based neural dialogue generator has been applied as an important component of practical chatbot services [16]. Thus, it is representative to be regarded as the object of study. However, due to the black-box setting of the dialogue model environment, we can manipulate other neural dialogue models such as HRED [14], VHRED [15], etc., indiscriminately through our proposed method. We leave the investigation into other choices as one future work.

#### 3.2 The Reverse Dialogue Generator Model

In order to find a corresponding input for a given targeted output, it is intuitive to formulate this problem as an optimization problem where we need to find an optimal input  $I$  so that the similarity between the output  $O = D(I)$  and the targeted output  $O^*$  are maximized. However, the standard gradient-based back-propagation approach is not applicable to this scenario. Because the input  $I$

consists of discrete tokens instead of continuous values, we cannot directly use gradient information to adjust it. If we make adjustments with regard to the word embeddings, this method will generate results that cannot be matched with any valid words in the word embedding space [29]. What’s more, in our black-box setting, we are not allowed to get the gradient information from  $D$ .

Therefore, in our work, we formulate this problem as a reinforcement learning problem. A generative agent  $G_\theta$ , which is called Reverse Dialogue Generator, is trained to craft an input  $I$  for a given targeted output  $O^*$ , with the aim to maximize the reward, that is, the similarity between the current output  $O$  and the targeted output  $O^*$ . The agent  $G_\theta$  treats the input generation process as a decision-making process, and it is trained through policy gradient with the guidance of the reward signals.

An RNN language model can be adopted as the agent  $G_\theta$ . However, we have to retrain it for each targeted output to obtain the desired input, which requires a great number of interactions with the  $D$  for every single targeted output. It is unachievable in reality when hackers try to manipulate a chatbot service. Instead, we adopt a classical Seq2Seq structure as the agent  $G_\theta$ , which takes a targeted output as input and crafts the desired input. It is trained offline, and when deployed in practical applications, it is able to generate the corresponding inputs for lots of targeted outputs automatically. The details of its training process are described in Section 4.2.

## 4 Training

In this section, we detail the training process of the dialogue model environment  $D$  and the Reverse Dialogue Generator model  $G_\theta$ .

### 4.1 The Dialogue Model Environment

The Seq2Seq dialogue model is trained on a large-scale human conversation data collected on Twitter in a supervised manner with the aim to minimize the negative log-likelihood of the ground truth sentence given inputs. Afterward, it is treated as the black-box environment and its parameters are not further updated. The details of the implementation can be found in Section 5.1.1.

### 4.2 RL Training of Reverse Dialogue Generator Model

In this subsection, we detail the RL training process of the proposed Reverse Dialogue Generator agent  $G_\theta$ . In summary, the agent (Reverse Dialogue Generator  $G_\theta$ ) interacts with the environment (the dialogue model  $D$ ). Given a state (the targeted output  $O^*$ ), the agent takes an action (generating an input  $I$ ) following a policy  $\pi_\theta$  (defined by the Seq2Seq model in the Reverse Dialogue Generator), and then receives a reward (the similarity between the targeted output  $O^*$  and current output  $O$ ) from the environment. Afterward, the policy  $\pi_\theta$  is updated to maximize the reward.

Next, we will introduce the environment, state, policy, action, and reward in detail.

#### 4.2.1 Environment

The environment is the black-box dialogue model  $D$ . When fed into an input  $I$ , it returns an output  $O = D(I)$ . The input and the output are two dialogue utterances, which consist of a sequence of words with variable length.

#### 4.2.2 State

A state is denoted as the targeted output  $O^*$ , that is, the input of the Reverse Dialogue Generator.

#### 4.2.3 Policy and Action

The policy  $\pi_\theta$  is defined by the Seq2Seq model in the agent  $G_\theta$  and its parameters. The Seq2Seq model forms a stochastic policy which assigns a probability to any possible input  $I = (\omega_1, \dots, \omega_T)$ :

$$\pi_\theta(I|O^*) = \prod_{t=1}^T P(\omega_t|\omega_1, \dots, \omega_{t-1}, O^*) \tag{1}$$

where  $T$  is the length of the input and  $\omega_t$  is the  $t$ -th word.

The action is defined as the input  $I$  to generate. When observing the state  $O^*$ , the agent generates an input based on the distribution predicted by  $\pi_\theta$  in Equation (1). In the training phase, the input is chosen by stochastic sampling. And in the test phase, the input can be chosen in a greedy manner or through beam search.

#### 4.2.4 Reward

Recall that our goal is to train the agent  $G_\theta$  to craft an input for the dialogue model  $D$  to return an output that is as similar with the targeted output  $O^*$  as possible. Let  $O$  be the actual output returned by the dialogue model  $D$  given a crafted input  $I$ . We directly use the similarity between the targeted output  $O^*$  and the current output  $O$  as the reward for the input  $I$  selected by  $\pi_\theta$ .

We adopt the embedding average metric to measure the similarity. This metric has been frequently used in many NLP domains, such as textual similarity tasks [36], since it’s able to measure the similarity of two sentences in semantic level rather than simply consider the amount of word-overlap. The embedding average approach first computes the sentence-level embedding  $\bar{e}_r$  of a sentence  $r$  by taking the average of the embeddings  $e_\omega$  of all the constituent words  $\omega$  in it:  $\bar{e}_r = \frac{\sum_{\omega \in r} e_\omega}{\|\sum_{\omega' \in r} e_{\omega'}\|_2}$ , and then the similarity between two sentences is defined as the cosine similarity of their corresponding sentence-level embeddings. Given a crafted input  $I$  and the targeted output  $O^*$ , we formally define the reward as the similarity between the current output  $O$  and the targeted output  $O^*$ :

$$R(I|O^*) = \text{Sim}(O^*, D(I)) = \cos(\bar{e}_O, \bar{e}_{O^*}) \quad (2)$$

#### 4.2.5 Optimization

With the reward function defined above, the objective function that the agent aims to maximize can be formulated as follows:

$$J(\theta) = \mathbb{E}_{I \sim \pi_\theta(I|O^*)} R(I|O^*) \quad (3)$$

The accurate value of  $J(\theta)$  in Equation (3) is very difficult to obtain in practice. Therefore, previous works have proposed many methods to estimate it and its gradient, which is then used to update the parameters  $\theta$  of the policy (i.e.,  $\pi_\theta$ ).

To optimize the objective in Equation (3), we apply the widely used REINFORCE algorithm [40], where Monte-Carlo sampling is applied to estimate  $\nabla_\theta J(\theta)$ . Specifically,

$$\begin{aligned} \nabla_\theta J(\theta) &= \sum_I R(I|O^*) \nabla \pi_\theta(I|O^*) \\ &= \sum_I R(I|O^*) \pi_\theta(I|O^*) \nabla \log \pi_\theta(I|O^*) \\ &= \mathbb{E}_{I \sim \pi_\theta(I|O^*)} [R(I|O^*) \pi_\theta(I|O^*) \nabla \log \pi_\theta(I|O^*)] \\ &\approx \frac{1}{N} \sum^N R(I|O^*) \pi_\theta(I|O^*) \nabla \log \pi_\theta(I|O^*) \end{aligned} \quad (4)$$

With the obtained gradient  $\nabla_\theta J(\theta)$ , the parameters of the policy network  $\pi_\theta$  can be updated as follows:

$$\theta := \theta + \alpha \nabla_\theta J(\theta) \quad (5)$$

where  $\alpha$  is the learning rate. Thus, the REINFORCE algorithm for updating the policy  $\pi_\theta$  can be summarized as: For each targeted output  $O^*$ , we first sample  $N$  inputs according to the distribution  $\pi_\theta(\cdot|O^*)$ . Then we estimate the rewards of the sampled inputs and calculate the gradient. Finally, we update the parameters of the policy network.

## 5 Experiments

We conduct extensive experiments to evaluate the effectiveness of the proposed Reverse Dialogue Generator. We measure the success rates of the proposed model in seeking out the desired inputs towards various targeted outputs and explore the performance of the model under various settings. In this section, we will go into details about the experimental settings and results.

## 5.1 Experimental Settings

### 5.1.1 The Dialogue Model

To ensure the reproducibility of the experiments, we directly adopt the well-trained Seq2Seq dialogue model released on the dialogue system research software platform ParlAI [41] as the dialogue model environment. The implementation of the dialogue model is detailed as follows. In the Seq2Seq model, both the encoder and the decoder are implemented by 3-layer LSTM networks with hidden states of size 1024. As the standard practice, the initial hidden state of the decoder is set as the same as the last hidden state of the encoder. The vocabulary size is 30,000. Pre-trained Glove word vectors [42] are used to initialize the word embeddings whose dimension is set as 300. The model had been trained through stochastic gradient descent (SGD) with a learning rate of 1.0 on 2.5 million Twitter single-turn dialogues. In the training process, the dropout rate and gradient clipping value are both set to be 0.1. It should be pointed out again that in the following experiments, this dialogue model is treated as a black-box which takes an input sentence and output a response sentence.

### 5.1.2 The Reverse Dialogue Generator

As described above, we adopt a Seq2Seq structure as the Reverse Dialogue Generator. Two 2-layer LSTM networks with the hidden size of 1,024 are applied as the encoder and the decoder respectively. The vocabulary size is set to be 60,000, and all the size of word embeddings is 300. The word embeddings are randomly initialized and fine-tuned during the training process. The last hidden state of the encoder is treated as the context vector which is used to initialize the hidden state of the decoder.

**Pre-training.** Before RL training, we first initialize the agent by pre-training it on output-input pairs in a reference corpus in a supervised learning manner. We build the reference corpus in this way: we first randomly collect 160K posts from Twitter, and then we feed them into the dialogue model to get the corresponding responses output by the model. Finally we obtain 160K output-input pairs as the reference corpus. Specifically, in the pre-training process, the model is optimized by the standard SGD algorithm with the initial learning rate of 20. At the end of each epoch, if the loss doesn't decrease in the validation set, the learning rate is reduced with a decay rate of 0.25. And the batch size is 16. In addition, to prevent overfitting issues, we apply the dropout with the rate of 0.1 and gradient clipping with clip-value being 0.25.

**RL training.** In the RL training process, all the parameters in the model are optimized as Equation (5) through Adam optimizer with an initial learning rate of 0.001. The 160K outputs in the reference corpus are used as the targeted outputs for RL training. The batch size, dropout rate, and gradient clipping value are set the same as those in the pre-training process. When calculating the rewards (i.e. the embedding similarities), we only consider the tokens out of a fixed stopword list which consists of common punctuations and special symbols. Pre-trained Google news word embeddings<sup>1</sup> are adopted to compute the similarities. And in order to accelerate the convergence of the model, we set all the rewards less than 0.5 to be 0, and the others remain the original values.

**Decoding.** In the test phase, the greedy method and beam search can be used for decoding. And we empirically find that if we first get  $N$  candidates with top- $N$  scores in beam search and then feed them into the dialogue model, treat the one with the highest reward as the crafted input, the performance improves significantly. In the experiments, we report the results of greedy decoding of the pre-trained model (Pre-trained Greedy), greedy decoding of the RL-trained model (RL Greedy), and beam search of the RL-trained model with  $N$  candidates (RL BeamSearch( $N$ )).

## 5.2 Experimental Results

We conduct experiments on two types of target output lists. To construct the **Generated** target output list, we first feed 10k human utterances (have no overlap with the reference corpus) from Twitter into the dialogue model to get 10k generated responses and then randomly sample 200 responses as the targets in length 1-3, 4-6 and 7-10, respectively. The **Real** target output list is obtained by randomly collecting sentences directly from Twitter. The number of target outputs in each length group is also 200.

---

<sup>1</sup><https://code.google.com/archive/p/word2vec/>

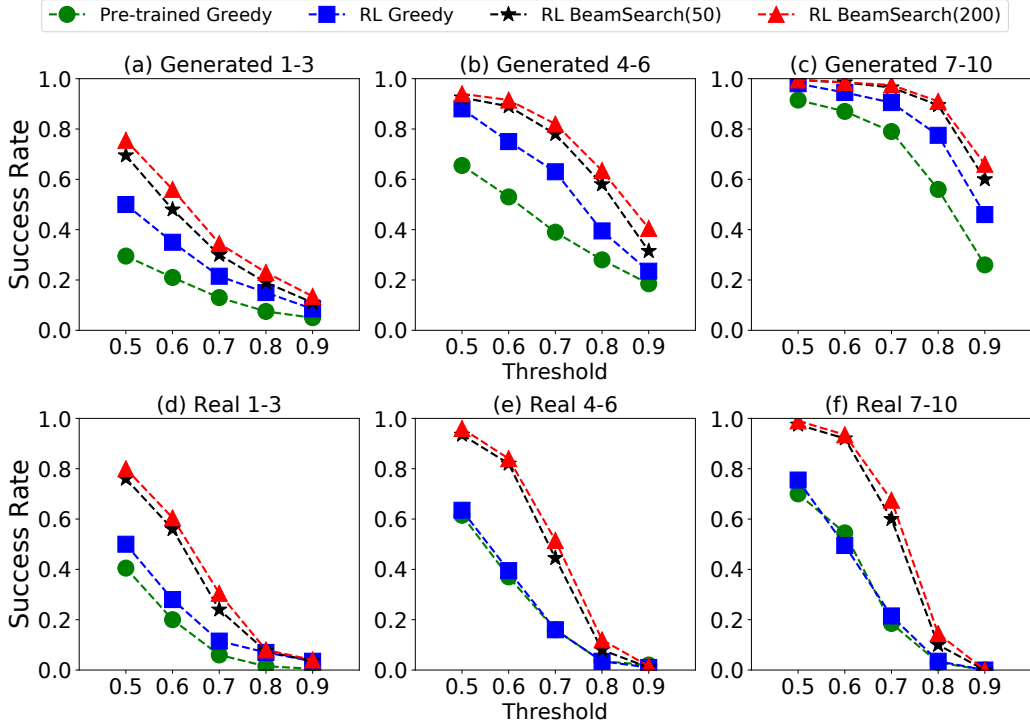


Figure 2: Success rates of the pre-trained model and the RL-trained model with different decoding methods. The upper row and the lower row show the results on the Generated target list and the Real target list with various lengths respectively.

Given a targeted output, when the proposed model finds an input which leads to an output whose similarity to the targeted one is above a preset threshold, we say the manipulation is successful. Figure 2 shows the success rates of our proposed model for manipulating the Twitter dialogue model in two experimental settings. The figures show how success rates vary with different thresholds. First of all, from the figures, we can see that for both the **Generated** and the **Real** target lists, RL-based model with beam search can achieve a success rate around 80% – 100% with a score 0.5. Especially for more than around 80% **Generated** targets with length greater than or equal to 4, we can find desired inputs that lead to a similarity score above 0.7. Second, we can see that in most of the scenarios, RL Greedy obviously outperforms Pre-trained Greedy, which demonstrates that using reward signals to fine-tune the policy network through RL significantly strengthens its capability for crafting desired inputs. Besides, compared with greedy decoding method, beam search with multiple candidates improves the success rates a lot. And compared with  $N = 50$ , the beam search methods with  $N = 200$  candidates improve the performance slightly, which means that interacting with the dialogue model for a reasonably small number of times can guarantee a considerable success rate. Another key observation is that the model performs significantly better on the **Generated** target list than on the **Real** target output list. Actually, the neural dialogue models suffer from the safe response problem [43]. Such models tend to offer generic responses to diverse inputs, which makes it hard for the model to provide a specific targeted response (often seen in real human conversations).

In order to further demonstrate the effectiveness of the proposed framework, for each **Real** targeted output, we feed its corresponding real inputs in the corpus into the dialogue model to check how similar the output responses and the target ones are. We calculate these similarity scores for each **Real** targeted output for the real inputs and the inputs found by the proposed model and report the average value in Table 1. According to the table, we observe that even inputting the real inputs, the similarity scores between the outputs and the target outputs are not high. Besides, with the crafted inputs from the proposed framework, these similarity scores are significantly improved. For example, for RL BeamSearch(200), the similarity is improved by 41.5%, 34.0% and 28.3% for the target outputs with length 1-3, 4-6 and 7-10, respectively.

Table 1: Average embedding similarity scores between the output and the target output in terms of **Real** target output list.

<b>Length</b>	1-3	4-6	7-10
Real Input	0.439	0.518	0.566
Pre-trained Greedy	0.446	0.529	0.559
RL Greedy	0.486	0.560	0.588
RL BeamSearch(50)	0.599	0.678	0.709
RL BeamSearch(200)	0.621	0.694	0.726

Table 2: Case Study. The first column shows the inputs found by RL BeamSearch(200) according to given target outputs. The middle column shows the target outputs and the outputs generated by the dialogue model to the inputs. The column on the right side shows the embedding similarity score between outputs and target outputs.

<b>Inputs</b>	<b>Responses</b>	<b>Similarity</b>
soo calm . you should be nervous .	<b>Target:</b> i ' m just trying to be a good person <b>Output:</b> i ' m not . i ' m just trying to be a better person .	0.952
guess she'll be invited .	<b>Target:</b> i ' m sure she ' ll be fine . <b>Output:</b> i ' m sure she ' ll be a good one .	0.946
neither ready pls	<b>Target:</b> i ' m not ready for this <b>Output:</b> i ' m not ready for this	1.0
how is nephew ?	<b>Target:</b> he ' s a good guy <b>Output:</b> he ' s good . he ' s a good guy .	0.982
you weren't invited .	<b>Target:</b> i was there <b>Output:</b> i was there .	1.0

Table 2 shows five examples in the manipulating experiments. The first three target outputs are from the **Generated** target output list, while the other two are from the **Real** target list. Given those target outputs, desired inputs are successfully crafted. Each of them leads to an output of the dialogue model similar or equal to the target one, evaluating by the embedding similarity measurement. Besides, unlike some related works [21, 37] where crafted text inputs are ungrammatical and meaningless, the inputs generated by our model are smooth and natural utterances. This is because the decoder of the Seq2Seq model in the Reverse Dialogue Generator serves as a language model, which guarantees the smoothness of the generated inputs.

## 6 Conclusion

Recently, dialogue systems are being integrated into our daily lives at a quite rapid pace. In the practical implementation of dialogue systems, neural dialogue models play an important role. However, recent concerns have risen for neural models across all domains as to whether they can be manipulated (in most cases, by crafted adversarial examples), which inspires us to examine the same problem of neural dialogue models. Our work reveals a dark side of such models that they are likely to be manipulated to "say what I want" – by finding well-designed inputs, we can induce the dialogue agent to provide desired outputs. We propose a reinforcement learning based **Reverse Dialogue Generator** which learns to craft such inputs automatically in the process of interacting with the black-box neural dialogue model. Extensive experiments on a representative neural dialogue model demonstrate the effectiveness of our proposed model and show that dialogue systems used in our daily lives can indeed be manipulated, which is a warning about the security of dialogue systems for both the research community and the industry.

For future works, we plan to extend the current framework to other sequence models. Besides, in this work, we examine the security problem of dialogue systems. In future works, we will also investigate concerns about the privacy of them, specifically, the possibility of dialogue systems to leak the sensitive information of users.



## References

- [1] Joseph Weizenbaum. ELIZA - a computer program for the study of natural language communication between man and machine. *Commun. ACM*, 9(1):36–45, 1966.
- [2] David Goddeau, Helen M. Meng, Joseph Polifroni, Stephanie Seneff, and Senis Busayapongchai. A form-based dialogue manager for spoken language applications. In *The 4th International Conference on Spoken Language Processing, Philadelphia, PA, USA, October 3-6, 1996*, 1996.
- [3] Bradford W. Mott, James C. Lester, and Karl Branting. Conversational agents. In *The Practical Handbook of Internet Computing*. 2004.
- [4] Rafael E. Banchs. Movie-dic: a movie dialogue corpus for research and development. In *The 50th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, July 8-14, 2012, Jeju Island, Korea - Volume 2: Short Papers*, pages 203–207, 2012.
- [5] David Ameixa, Luísa Coheur, Pedro Fialho, and Paulo Quaresma. Luke, I am your father: Dealing with out-of-domain requests by using movies subtitles. In *Intelligent Virtual Agents - 14th International Conference, IVA 2014, Boston, MA, USA, August 27-29, 2014. Proceedings*, pages 13–21, 2014.
- [6] Zhengdong Lu and Hang Li. A deep architecture for matching short texts. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 1367–1375, 2013.
- [7] Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen. Convolutional neural network architectures for matching natural language sentences. *CoRR*, abs/1503.03244, 2015.
- [8] Yiping Song, Rui Yan, Xiang Li, Dongyan Zhao, and Ming Zhang. Two are better than one: An ensemble of retrieval- and generation-based dialog systems. *CoRR*, abs/1610.07149, 2016.
- [9] Jianfeng Gao, Michel Galley, and Lihong Li. Neural approaches to conversational AI. *Foundations and Trends in Information Retrieval*, 13(2-3):127–298, 2019.
- [10] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 3104–3112, 2014.
- [11] Alessandro Sordani, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. A neural network approach to context-sensitive generation of conversational responses. In *NAACL HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015*, pages 196–205, 2015.
- [12] Oriol Vinyals and Quoc V. Le. A neural conversational model. *CoRR*, abs/1506.05869, 2015.
- [13] Lifeng Shang, Zhengdong Lu, and Hang Li. Neural responding machine for short-text conversation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 1: Long Papers*, pages 1577–1586, 2015.
- [14] Iulian Vlad Serban, Alessandro Sordani, Yoshua Bengio, Aaron C. Courville, and Joelle Pineau. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA.*, pages 3776–3784, 2016.
- [15] Iulian Vlad Serban, Alessandro Sordani, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron C. Courville, and Yoshua Bengio. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 3295–3301, 2017.

- [16] Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. The design and implementation of xiaoice, an empathetic social chatbot. *CoRR*, abs/1812.08989, 2018.
- [17] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- [18] Xiaoyong Yuan, Pan He, Qile Zhu, Rajendra Rana Bhat, and Xiaolin Li. Adversarial examples: Attacks and defenses for deep learning. *CoRR*, abs/1712.07107, 2017.
- [19] Marty J. Wolf, Keith W. Miller, and Frances S. Grodzinsky. Why we should have seen that coming: comments on microsoft’s tay "experiment," and wider implications. *SIGCAS Computers and Society*, 47(3):54–64, 2017.
- [20] Rob Price. Microsoft is deleting its ai chatbot’s incredibly racist tweets. *Business Insider*, 2016.
- [21] Tianxing He and James Glass. Detecting egregious responses in neural sequence-to-sequence models. *CoRR*, abs/1809.04113, 2018.
- [22] Hongge Chen, Huan Zhang, Pin-Yu Chen, Jinfeng Yi, and Cho-Jui Hsieh. Attacking visual language grounding with adversarial examples: A case study on neural image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 2587–2597, 2018.
- [23] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K. Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016*, pages 1528–1540, 2016.
- [24] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan L. Yuille. Adversarial examples for semantic segmentation and object detection. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1378–1387, 2017.
- [25] Jernej Kos, Ian Fischer, and Dawn Song. Adversarial examples for generative models. In *2018 IEEE Security and Privacy Workshops, SP Workshops 2018, San Francisco, CA, USA, May 24, 2018*, pages 36–42, 2018.
- [26] Pedro Tabacof, Julia Tavares, and Eduardo Valle. Adversarial images for variational autoencoders. *CoRR*, abs/1612.00155, 2016.
- [27] Sandy Huang, Nicolas Papernot, Ian J. Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*, 2017.
- [28] Jernej Kos and Dawn Song. Delving into adversarial attacks on deep policies. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*, 2017.
- [29] Wei Emma Zhang, Quan Z. Sheng, and Ahoud Abdulrahmn F. Alhazmi. Generating textual adversarial examples for deep learning models: A survey. *CoRR*, abs/1901.06796, 2019.
- [30] Qi Lei, Lingfei Wu, Pin-Yu Chen, Alexandros G. Dimakis, Inderjit S. Dhillon, and Michael Witbrock. Discrete attacks and submodular optimization with applications to text classification. *CoRR*, abs/1812.00151, 2018.
- [31] Bin Liang, Hongcheng Li, Miaoqiang Su, Pan Bian, Xirong Li, and Wenchang Shi. Deep text classification can be fooled. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 4208–4215, 2018.

- [32] Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. Adversarial example generation with syntactically controlled paraphrase networks. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 1875–1885, 2018.
- [33] Motoki Sato, Jun Suzuki, Hiroyuki Shindo, and Yuji Matsumoto. Interpretable adversarial perturbation in input embedding space for text. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden.*, pages 4323–4330, 2018.
- [34] Yonatan Belinkov and Yonatan Bisk. Synthetic and natural noise both break neural machine translation. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- [35] Alvin Chan, Lei Ma, Felix Juefei-Xu, Xiaofei Xie, Yang Liu, and Yew Soon Ong. Metamorphic relation based adversarial attacks on differentiable neural computer. *CoRR*, abs/1809.02444, 2018.
- [36] John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. Towards universal paraphrastic sentence embeddings. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- [37] Minhao Cheng, Jinfeng Yi, Huan Zhang, Pin-Yu Chen, and Cho-Jui Hsieh. Seq2sick: Evaluating the robustness of sequence-to-sequence models with adversarial examples. *CoRR*, abs/1803.01128, 2018.
- [38] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [39] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. In *Proceedings of SSST@EMNLP 2014, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation, Doha, Qatar, 25 October 2014*, pages 103–111, 2014.
- [40] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [41] Alexander H. Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, Devi Parikh, and Jason Weston. Parlai: A dialog research software platform. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017 - System Demonstrations*, pages 79–84, 2017.
- [42] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [43] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting objective function for neural conversation models. In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 110–119, 2016.