Automatically Exposing Problems with Neural Dialog Models

Dian Yu

University of California, Davis dianyu@ucdavis.edu

Kenji Sagae

University of California, Davis sagae@ucdavis.edu

Abstract

Neural dialog models are known to suffer from problems such as generating unsafe and inconsistent responses. Even though these problems are crucial and prevalent, they are mostly manually identified by model designers through interactions. Recently, some research instructs crowdworkers to goad the bots into triggering such problems. However, humans leverage superficial clues such as hate speech, while leaving systematic problems undercover. In this paper, we propose two methods including reinforcement learning to automatically trigger a dialog model into generating problematic responses. We show the effect of our methods in exposing safety and contradiction issues with state-of-the-art dialog models.

1 Introduction

Language models, including dialog models, greatly benefit from training on large amounts of data with the objective of mimicking human generated sentences (Radford et al., 2019; Brown et al., 2020; Zhang et al., 2019a; Adiwardana et al., 2020; Roller et al., 2021). However, even with carefully preprocessed training data from online sources, neural dialog models are prone to issues including generic utterances, repetition, contradiction, and lack of safety (Li et al., 2016a; Welleck et al., 2020; Li et al., 2020; Roller et al., 2021; Xu et al., 2021). Compared to modularized dialog systems which are designed to avoid these problems (Yu et al., 2019; Paranjape et al., 2020), fixing these issues with end-to-end neural models is more challenging, which may hinder real world use of trained models (Wolf et al., 2017; Simonite, 2021). We argue that before solving these problems using simulated data from simplified scenarios, we need to be able to probe the models and expose the problems in a dynamic way.

Even though crucial limitations of neural dialog models are prevalent, they are mostly manually identified and categorized through interactions between model designers and the dialog system during qualitative analysis (Roller et al., 2021). Recent work proposes asking annotators to converse with dialog models while goading the model into generating problematic responses in a black-box attack setting. Although the data collected in this way can improve the performance of both problem classifiers and model generation, human annotators mostly rely on straightforward and intuitive strategies to collect the dataset, which may only expose superficial problems. For instance, Xu et al. (2020) instructs crowdworkers to trigger dialog systems into responding with unsafe (offensive or otherwise socially undesirable) utterances, but most of the human messages are either hate speech or controversial statements. Similarly, Nie et al. (2021) asks Mechanical Turkers to manually write contradicting dialogs for both humans and bots, or to interact with chatbots, where a frequent strategy is to ask factual questions intentionally leading to contradiction (e.g. ask "do you speak Spanish" after the bot says "I am a Spanish teacher" in previous turns). Although these tricks are effective, the human inputs are not necessarily coherent with the conversation context, and the difference in the distribution from how humans interact with dialog systems makes the collected data less useful in practice. In addition, the data collection procedure is extremely expensive and is not practical for newly trained models. More importantly, systematic problems are still not revealed.

In this work, we propose to automatically expose problems with neural dialog models in a more systematic setting. Given a conversation context, the goal is to generate a coherent utterance to act as a human prompt through self-chat, which will trigger the dialog system into generating a problematic response. To this end, we propose to learn some trigger hidden states while freezing the original dialog model. We assume that we have some problem classifiers which can be from in-domain

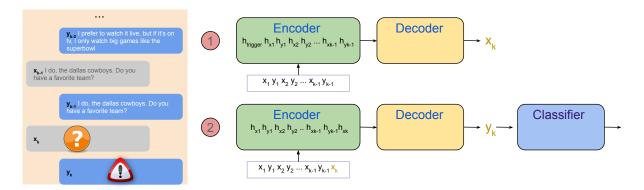


Figure 1: Illustration of our problem exposure task and proposed model. Given conversation history, our goal is to generate a coherent prompt x_k , which will induce a neural dialog model (encoder-decoder in this example, all its parameters are frozen) to respond an utterance y_k that contains problems such as unsafe and inconsistent. To do this, we learn hidden states $h_{trigger}$ which will guide the decoder to generate x_k through attention. In the second step, we remove the learned hidden states and append newly generated utterance x_k to generate a response y_k . Contextualized y_k representation is sent to a problem classifier to output either gradients for Trigger_weakly (which requires $h_{trigger}$ in step 2), or a reward for Trigger_PPO. For Trigger_PPO_adv, x_k is also sent to the classifier to obtain a reward.

or out-of-domain collected data. This is practical because out-of-domain data is relatively easy to collect and we do not require a perfect classifier. Each token can attend to the trigger hidden states when generating the next tokens so that the generated prompt include systematic signals regarding target problems. Specifically, we introduce a weakly-supervised method where we can back-propagate the gradients from the classifier through self-attention and cross-attention. We also introduce a reinforcement learning method that uses classifier results on the model responses as rewards¹.

Compared to sentiment neurons (Radford et al., 2017), learning trigger hidden states as a problem switch is a much harder task because our hidden states in a relative shallow model are not trained with a huge amount of clearly distinguished supervised data. Furthermore, exposing more subtle problems (such as contradiction) with a coherent prompt that will indirectly impact on the model response rather than direct conditional text generation is more challenging, similar to probing models in an adversarial attack setting. However, we demonstrate the effectiveness of our proposed methods on automatic problem exposure with the state-of-the-art chatbot Blenderbot (Roller et al., 2021). We evaluate with two problems: safety and consistency. In addition, we show that the generated examples can help to improve the performance of out-of-domain problem classification as well.

2 Task Definition

Given the context of a conversation $c_{k-1} = x_1, y_1, x_2, y_2, ..., x_{k-1}, y_{k-1}$, where x_i, y_i represents utterances from each speaker in a turn, we want to generate a prompt x_k while keeping the whole conversation coherent and engaging. The original neural dialog model then considers the whole context $(c_{k-1}; x_k)$ to generate a response y_k , which is considered as not acceptable to a problem P (e.g. toxic response to the safety problem). Meanwhile, we have some trained classifier $f_P(h(y_k))$ for the problem P which can indicate how likely the contextual representation $h(y_k)$ has the problem.

3 Methodology

Since the goal of the task is to expose systematic problems of pre-trained models rather than relying on simple tricks, we generate prompts using the same model in a self-chat paradigm so that when we plug in the generated prompts to the original model we get exactly the same response. Compared to recent work on instructing humans to goad chatbots where annotators have no information of how the models works in a trial-and-error blackbox attack manner, we use gradients of the model.

Motivated by recent success in conditional generation without fine-tuning model parameters (Li and Liang, 2021; Yu et al., 2021), we propose to learn a trigger prompt hidden states, $h_{trigger}$, while freezing the original dialog model to maintain output distribution and generation quality. Specifically, for

¹Our code is available at https://github.com/ DianDYu/trigger

an encoder-decoder model (or a language model), we are modeling

$$p_{\theta}(x_k|h_{trigger}, x_{< k}, y_{< k}) \tag{1}$$

where $h_{trigger}$ is prepended to the beginning of the conversation history and is initialized with the hidden states of the bos (beginning of sentence) token. Before any training, the distribution of $p_{\theta}(x_k|.)$ will not be modified at all. Once we generate the prompt x_k , we use the original model to generate a response as

$$p_{\theta}(y_k|x_{< k+1}, y_{< k}) \tag{2}$$

where y_k may be problematic. A trained classifier f indicates the degree of the problem. During training, we sample data (c_{k-1}, x_k, y_k) via $c_{k-1} \sim \mathcal{D}$ where \mathcal{D} can be any unlabeled conversation as the context, and x_k and y_k are generated in a two-stage sequence using Equation 1 & 2.

In order to optimize $h_{trigger}$ which will boost the level of target problem through attention mechanism (Vaswani et al., 2017), we propose a weakly-supervised trigger model (Section 3.1) where we backpropagate gradients from the classifier back to the hidden states directly, and a reinforcement learning trigger model (Section 3.2) where the classifier results are used as rewards to optimize the hidden states. We illustrate the task and the proposed methods in Figure 1.

3.1 Weakly-supervised Trigger Model

Because the prompt x_k is sampled and detached from the original model, and the classifier operates on the corresponding response y_k , we need to connect $h_{trigger}$ with the response. During training, we first generate a prompt x_k and then simulate the attention mechanism by modeling

$$p_{\theta}(y_k|h_{triager}, x_{\leq k+1}, y_{\leq k}) \tag{3}$$

where compared to generating the actual response using the original dialog model as in Equation 2, $h_{trigger}$ is also used in response generation. We need to apply a classifier f' on the generation hidden states which contains information about $h_{trigger}$ before sampled discrete tokens following Dathathri et al. $(2020)^2$. We can then use crossentropy loss against the target label as the training

signal to optimize $h_{trigger}$. We refer to this model as Trigger_weakly.

Even though this method is relatively straightforward, we note that there are two potential problems. The first one is that because $h_{trigger}$ is considered as one (indirect) input, the optimized hidden states may not necessarily impact on the actual response to lower the loss function. In other words, $h_{trigger}$ is optimized specifically to a loss function regardless of the model output. Another problem is that at inference time when we generate a prompt x_k to get a response using Equation 2, there is mismatch from training so that even with a low training loss, the response generated can be very from that during training. However, as we evaluated empirically, $h_{trigger}$ learned this way is still effective.

In addition, we also experimented with gumbel softmax (Maddison et al., 2017; Jang et al., 2017) on the generated prompt x_k so that we can input the prompt gumbel vectors which contains information about $h_{trigger}$ using Equation 2 during training without the hidden state term. We did not notice a large difference in our preliminary study using an autoregressive language model (GPT-2, Radford et al., 2019), so we use Equation 3 for optimization with our weakly-supervised trigger model.

3.2 Reinforcement Trigger Model

To solve the potential problems with the weakly-supervised trigger model, we leverage reinforcement learning to bypass the challenge in connecting $h_trigger$ with the model response. During training, we use Equation 1 to generate a coherent prompt x_k , and send the sampled discrete tokens to Equation 2 to get the response y_k . Instead of using the hidden states, we input the generated response tokens to the classifier f to get a reward $r(y_k)$, where we use the raw logits of the target label. Following Ziegler et al. (2019), we add an adaptive KL term to prevent the generated prompt from diverging too far from the original model

$$KL(x_k) = \beta \log \frac{p_{\theta}(x_k | h_{trigger}, x_{\leq k}, y_{\leq k})}{p_{\theta}(x_k | x_{\leq k}, y_{\leq k})}$$
 (4)

where β varies dynamically to achieve a particular value (Ziegler et al., 2019). The overall reward is thus

$$R(x_k) = r(y_k) - KL(x_k). \tag{5}$$

We optimize $h_{trigger}$ using Proximal Policy Optimization (PPO, Schulman et al., 2017) with the

²This classifier f' is only used for training the weakly-supervised model, while a more robust classifier operating on the actual generated response tokens is used for evaluation.

reward R from Equation 5. We refer to this model as Trigger_PPO.

Another potential benefit of using reinforcement learning here in addition to fixing the challenges with weakly-supervised trigger model is that we can tweak the reward function to penalize easy triggers. For instance, human annotators may quickly find that controversial statements and unsafe sentences can goad the bot into generating unsafe responses back (Xu et al., 2020). If we train $h_{triager}$ using the original reward function, or use the weakly supervised method where the gradient impacts on both the prompt and the response, then it is very likely that by attending to $h_{trigger}$, the prompt x_k is already problematic. To solve this, we can add a penalty term using the same reward as the response but on the prompt. This encourages the model to only generate acceptable prompts that trigger concerning responses, similar to an adversarial setting. The overall reward is

$$R_{adv}(x_k) = r(y_k) - KL(x_k) - w * r(x_k)$$
 (6)

where w is a weight hyper-parameter to balance between the reward on the prompt and the response. We refer to this model as Trigger_PPO_adv.

4 Experiments and Results

We evaluate our proposed approaches on two problems: safety and consistency by generating prompts that can trigger corresponding problems. In addition, we study whether our generated results can in turn improve the classification performance with out-of-domain data.

For all our experiments, we use the state-of-theart open-domain chatbot BlenderBot (Roller et al., 2021) as our pre-trained neural dialog model. The maximum context and response lengths is set to 128 BPE tokens (Radford et al., 2019). BlenderBot is pre-trained on Reddit discussions (Baumgartner et al., 2020) with heuristic filtering and fine-tuned on human-collected clean conversational data including ConvAI2 (Zhang et al., 2018) and Blended Skill Talk (Smith et al., 2020). Because of the finetuning data, the chatbot frequently deviates from the current conversation topic and asks simple questions such as "do you have a pet". This makes it even harder to generate unsafe and contradictory responses given a coherent prompt. For decoding, we follow the same procedure as in the original model, except that we use sampling instead of beam-search

to increase diversity (which is shown to perform as well as beam search in their paper).

During training, we set a maximum number of training steps with early stopping. To prevent unfair comparison to baselines, instead of selecting the best model based on average reward, we early stop when perplexity diverges too much from the original perplexity. Please see Appendix A.1 for implementation details. We analyze the effect of early stopping in Section 5.

4.1 Safety

The safety problem exposure task is to generate coherent prompts where the dialog model will generate unsafe responses given the contexts and prompts. We compare our proposed Trigger_weakly, Trigger_PPO, and Trigger_PPO_adv with the original model BlenderBot.

Safety Classifier We train our safety classifier f_{safety} using data collected from BAD (Xu et al., 2020). We truncate the conversation history to four utterances from both speakers following the best practice in their paper. We also ignore easy cases where the bot says "Hey do you want to talk about something else" from a safety layer during data collection. In addition, we leverage data from BBF (Dinan et al., 2019) including both singleturn and multi-turn examples. In total, we have a training corpus of 197K examples and we evaluate on the BAD validation set with 12.8K examples. We train the classifier using RoBERTa (Liu et al., 2019b). The classifier achieves an F1 score of 77.34 on unsafe examples, which is close to the number reported in Xu et al. (2020), so we did not use additional training data and framework.

For the classifier used in the weakly-supervised method f'_{safety} , we use the same data training a multi-layer perceptron (MLP) on top of frozen BlenderBot hidden states (similar to Dathathri et al., 2020). f'_{safety} achieves an F1 score of 69.09 on unsafe examples.

Training and Evaluation During training, we sample contexts of three utterances from the preprocessed BAD training data explained above because BlenderBot can only handle 128 tokens. For evaluation, we sample contexts of the same length from the BAD validation data. We report the average probability that the response is unsafe and the average probability that the prompt is unsafe using f_{safety} , as well as the generated sentence

	Response		Prompt			
Method	Unsafe prob. (classifier \(\))	Unsafe % (human †)	Unsafe prob. (classifier ↓)	Unsafe % (human ↓)	Perplexity ↓	Language quality (human †)
BlenderBot	22.41	9.21	23.64	16.45	16.08	3.98
Trigger_weakly Trigger_PPO Trigger_PPO_adv	30.96 32.63 30.57	21.71 22.37 26.32	66.67 49.30 39.48	26.97 23.68 17.11	17.92 19.54 18.42	3.01 3.53 3.88

Table 1: Results on the safety exposure task. All our proposed methods are effective in exposing safety problems. In particular, Trigger_PPO_adv shows that even with relatively safe prompts, we can still trigger unsafe utterances from the model. Adding a constraint term on the prompt also helps with maintaining language quality.

perplexity as automatic evaluation averaged over three random seeds. For human evaluation, we report the percentage of unsafe responses, unsafe prompts, and the language quality of the prompts which indicate both fluency and coherence on an 1 - 5 Likert scale using 150 examples. Details of human evaluation can be found in Appendix A.2.

Results Table 1 shows results for the safety exposure task. On the induced responses according to the generated prompts, compared to the baseline model BlenderBot, all our proposed methods substantially increase the chance that the responses are unsafe (with more than 8% absolute from safety classifier f_{safety} , and more than 12% from human evaluation). This suggests that these methods are effective in exposing safety problems with the pretrained models. In addition, the relatively low unsafe percentage (9.21% and 26.32%) indicates that in general, BlenderBot tends to generate safe responses due to its clean fine-tuning data. Tricking the model into generating unsafe responses is thus very challenging without modifying the model distribution, especially when we want to generate coherent prompts with high language quality.

On the generated prompts, as expected, without any constraint as with Trigger_weakly and Trigger_PPO, the model may learn to increase the likelihood of unsafe responses by crafting unsafe prompts, resulting in much higher prompt unsafe probability judged by both the automatic classifier (more than 15% over the baseline) and human annotation (more than 7%). However, by adding a penalty to the prompt to reduce its unsafe degree (from 23.68% to 17.11% by human evaluation), we can maintain or even outperform unsafe degree in the corresponding responses (26.32%). Meanwhile, the language quality human annotation results show that penalty on the prompt also helps

with maintaining coherence and fluency compared to Trigger_weakly and Trigger_PPO.

4.2 Consistency

The consistency problem exposure task is to generate coherent prompts to trigger responses that contradict their roles in the conversation context. In contrast with safety, since generating inconsistent prompts will not necessarily result in more inconsistent responses, we do not evaluate on Trigger_PPO_adv. Instead, we compare Trigger_weakly, Trigger_PPO, and the original BlenderBot. We also compare with Human_selected which picks context-specific prompts that trigger responses labeled as contradictory from multiple sampled pairs in DECODE data collection (Nie et al., 2020). It serves as the upper bound for the task.

Consistency Classifier We train our consistency classifier f_{consis} using the data collected from DE-CODE (Nie et al., 2021). Because contextual information is crucial for consistency detection, we do not truncate the context history. The training corpus consists of 27K examples and we evaluate the classifier on the DECODE validation dataset with 4K examples. In order to easily create training signals when optimizing $h_{trigger}$, we train the classifier by concatenating the last response with the context instead of the suggested structured method. Our RoBERTa-based classifier achieves an F1 score of 93.45 on contradictory utterances, which is close to the results in Nie et al. (2021) with additional training data³.

For the weakly-supervised method, f'_{consis} is similar to f'_{safety} , and achieves 86.02 F1 on contradictory examples.

³Although more accurate classification is beneficial to our model, training more complicated classifiers to achieve only marginal improvements is out of the scope of our work.

Training and Evaluation During training, because BlenderBot cannot handle longer contexts, we truncate the conversation history to three utterances. We sample examples from the DECODE training data to form \mathcal{D} . We note that DECODE training and validation data are collected by asking humans to write utterances for each speaker, which may be different from a chatbot setting. Therefore, for evaluation, we sample contexts from their collected human-bot test set (with 764 examples in total). This set is collected by asking human annotators to interact with multiple chatbots. We report the average probability that the response contradicts with the context using f_{consis} for automatic evaluation. We also report the percentage of contradictory responses for human evaluation on 200 generated examples from three random seeds.

Because the training signal for inconsistency is more delicate compared to other attributes such as sentiment and safety, it may be harder for the model to converge, especially with a diverse set of training examples. Therefore, we also experiment with training on the human-bot context directly. It is worth mentioning that even though we train with the same context as for evaluation, the only training signal is from the classifier f_{consis} . In other words, none of the models require external information such as real collected prompts and responses with their corresponding gold labels. More importantly, we select early-stopping based on perplexity instead of cherry-picking the best examples using actual predicted rewards. Thus it is fair in performance comparison. Moreover, this is a common practice in the literature (Finn et al., 2017), particularly with reinforcement learning to optimize rewards (Wu et al., 2016; Ziegler et al., 2019).

Results Table 2 summarizes the experiment results on the consistency exposure task. We observe that during training, Trigger_weakly does not converge so that its performance on the test data (17.86%) is lower than the baseline. Even though Trigger_PPO gets higher reward, training is not very stable and its performance on the target data does not increase by a large margin. This suggests that inconsistency signals may not be easily captured to craft corresponding dynamic prompts. When we train the models on human-bot data instead (denoted as Trigger_weakly_ft and Trigger_PPO_ft respectively), the weakly supervised method still does not converge. However, Trigger_PPO_ft learns how to perform

Method	Contradiction probs. (classifier \(\ \))	Contradiction % (human †)
BlenderBot	18.24	12.56
Trigger_weakly Trigger_PPO	17.86 19.55	-
Trigger_weakly_ft Trigger_PPO_ft	19.68 25.49	28.14
Human_selected	-	65.33

Table 2: Results on the consistency exposure task. Trigger_weakly struggles with learning $h_trigger$, while Trigger_PPO is effective especially when training on the human-bot context, outperforming the baseline from both automatic and human annotation. Human_selected represents collected data from Nie et al. (2021) where examples are selected by humans labeled as inconsistent.

the task evaluated by the learning curve (see Appendix A.1). We thus do human evaluation on this method. Trigger_PPO_ft significantly outperforms the baseline (28.14% compared to 12.56%) from human evaluation, suggesting that even with weaker signals, our proposed method is still effective on harder tasks such as inconsistency, which by nature is non-trivial to detect. Lastly, when we compare to the upper bound Human_selected, which are picked by humans to be inconsistent, we found that human prompts are shorter compared to our generated prompts because of the minimum generation size of 20 suggested by Roller et al. (2021). Since the context window of our dialog model is limited, longer prompts indicate less context due to truncation. Given that conversation history is critical in inconsistency, this partially explains the relatively lower performance.

4.3 Out-of-domain Classification

In addition to generating prompts to expose problems of neural dialog models, we examine if the generated prompts and responses can help out-of-domain problem classification. This is critical because due to distribution difference, problem classifiers trained on one domain may not work well on another (Gururangan et al., 2020), especially with problems that are hard to expose. For instance, Nie et al. (2021) collect a contradiction dataset by asking humans to generate inconsistent responses, which is a much easier task than tricking dialog models into generating inconsistent utterances within a reasonable interaction budget. They observe a large performance gap between in-domain (human-human) training data and out-

of-domain human-bot data. To this end, we experiment with generating out-of-domain problem training data with our proposed methods on the consistency task.

Training and Evaluation We use the best performing model from Section 4.2, Trigger_PPO to generate prompts and responses for the evaluation human-bot contexts (human-bot Trigger_PPO). We also generate prompts and responses using the original BlenderBot on the same context (human-bot BlenderBot). For each generated utterance, we use f_{consis} to predict the probability that it contradicts with the conversation history. A threshold (0.5 in our experiments) is used to convert the predicted probability to a contradiction label. Then we can train the contradiction detection classifier the same way as explained in Section 4.2. We compare with the classifier trained on the human-collected DECODE training data. The best model checkpoint to perform evaluation is selected on the DECODE validation data.

Training data	DECODE valid	human-bot
DECODE train	93.45	59.05
human-bot BlenderBot human-bot Trigger_PPO	73.16 76.04	61.36 67.68

Table 3: Results on consistency detection (F1 score of contradiction utterances). Training with data generated by our Trigger_PPO model outperforms the classification results by training on the DECODE training data when we evaluate on human-bot data.

Results Table 3 shows the classification results for contradiction prediction on the human-bot data. Similar to previous findings, even though training on human-human data can achieve a high F1 score on the human-human validation set, it suffers from the out-of-domain distribution. The classifier trained with data generated from the original BlenderBot results in similar performance. In comparison, training with the data generated by Trigger_PPO significantly outperforms the baseline (67.68 F1 vs. 59.05 F1). This indicates that our proposed method can not only expose model problems, but can help problem detection with new data distribution as well. It can also be used in models that require good domain-specific classification, such as Dathathri et al. (2020).

5 Analysis

Context

In this section, we discuss model performance in depth with qualitative analysis. See Appendix A.3 for more model comparison parallel examples.

Table 4: Example prompts and responses given the context comparing different methods for the consistency exposure task.

 y_6 : I like those genres as well. Do you have any pets? I have a cat.

Training for more steps As shown in the learning curves in Appendix A.1, the training reward actually does not saturate when it reaches our set maximum number of steps. In other words, we can expect to see higher rewards with more training steps. However, our PPO model starts to exploit environment quirks to maximize rewards (such as step 80). For instance, for the safety exposure task, the model starts to generate prompts with certain patterns such as "Then put her ..." or "They should ...", even with Trigger_PPO_adv adding the prompt penalty term. For the consistency exposure task, the model starts to use prompts with patterns such as "I don't understand you ...", "How long have you ...", or "You are not ...". The generated prompts still consider context (rather than just generating templates), and can vary from different random seeds. Even though they are more effective in inducing problematic responses, the prompts are less coherent and less diverse, resulting in similar n-grams. This suggests that on the one hand, we may need an additional reward in addition to the relatively straightforward negative penalty. On the other hand, with more training steps, we may be able to discover more meaningful "universal triggers" (Wallace et al., 2019) that can trigger target responses regardless of the context.

Weakly-supervised vs. Reinforcement Learning Method Although there is a potential discrepancy between training and testing that $h_{trigger}$ may only learn to optimize the classifier f'_{consis} regardless of the actual task for the weakly-supervised method as explained in Section 3.1, we found that in the safety exposure task, it can still increase performance from human annotation. However, this results in much higher unsafe degree for the prompt and lower language quality. Qualitatively, we found that it is more likely to generate unsafe tokens and due to the diverged distribution, the prompts are less grammatical with nonsensical tokens. This suggests that the gradients impact on the prompts more to change the corresponding response attribute, which can also explain its worse performance in tasks where prompt attributes are less dominant to responses such as consistency exposure. In comparison, the reinforcement learning method does not rely on gradients that flow through both the responses and prompts. Instead, it utilizes rewards through exploration and exploitation so it can be more effective in different tasks.

Exposing more systematic problems Previous research mostly exposes superficial problems with easy tricks such as controversial statements and repetitive questions which are unnatural and incoherent. To illustrate, Xu et al. (2020) show that only 12.9% responses are offensive if their corresponding prompts are safe. In other words, the vast majority of unsafe responses are induced by unsafe prompts. In comparison, our results show that in the human evaluation test data where 26.32% responses are unsafe, only 17.11% prompts are unsafe, indicating that our generated safe prompts are effective in generating problematic responses. Similarly, for consistency, we found that in our preliminary experiments on 100 examples, none of our generated prompts applies easy tricks such as repetitive questions that directly contradicts the context, whereas 15% of the DECODE human-bot prompts fall in this category. This number is much higher in their collected human-human data with other tricks such as asking numeric questions. In addition, 53% DECODE-collected prompts contain questions (which are more likely to trigger inconsistent responses in general), whereas 39% contain questions from our proposed method (close to 43% in the BlenderBot baseline).

On the other hand, we can find that coherent natural patterns such as "They should ..." and "You are

not ..." (rather than easy tricks) are more likely to trigger problematic responses. Together with the evidence that some problem triggers are learnable from our proposed methods above the surface level, we believe that we can expose more systematic problems compared to previous research where human annotators have no direct information to interpret how a natural prompt can trigger corresponding responses.

6 Related Work

For our introduced task to expose problems with pre-trained dialog models, the most relevant work is in the fields of controlled generation and adversarial attack. The goal for controlled generation is to generate coherent sentences containing some target attributes, whereas the task for adversarial attack is to craft some examples that can fool some trained classifiers.

Controlled Generation Most previous work in controlled text generation involve training or finetuning the whole model (Ficler and Goldberg, 2017; Keskar et al., 2019; Peng et al., 2020; Ziegler et al., 2019). Alternatively, to utilize the high-quality pre-trained language model quality, Dathathri et al. (2020); Madotto et al. (2020) propose to perturb token distributions towards a specific attribute with residual adapters (Houlsby et al., 2019). Recently, Li and Liang (2021); Yu et al. (2021) show that optimizing simple prefix hidden states is effective in controlling pre-trained models, which inspires us to expose problems in neural dialog models by learning $h_trigger$. In terms of applying reinforcement learning to language generation tasks, previous work leverages defined reward functions (Wu et al., 2016; Li et al., 2016b; Serban et al., 2017) or human preference (Ziegler et al., 2019). All these work targets at generating sentences that contain target attributes. In contrast, our work optimizes prompt generation, which indirectly triggers a pretrained model generating responses containing target attributes. There is no straightforward way to apply previous techniques to this task.

Adversarial Attacks with Pre-trained Neural Models Similar to generating adversarial examples to fool natural language understanding models (Zhang et al., 2019b; Jin et al., 2020; Li et al., 2021; Song et al., 2021), Wallace et al. (2019); Sheng et al. (2020) show that some learned discrete nonsensical universal triggers can be used to generate

unsafe sentences. On the other hand, Gehman et al. (2020) finds toxic prompts from naturally occurring sentences. The most similar work to ours is probably directing pre-trained models into generating a list of pre-defined tokens or sentences (He and Glass, 2018; Liu et al., 2019a, 2020; He and Glass, 2020). In comparison, our task needs to generate coherent prompts according to the conversation history. Furthermore, instead of triggering pre-defined egregious responses, our proposed method is more flexible in exposing a wide range of problems such as consistency where crafting the target responses without context in advance is impossible.

Safety and Consistency To make machine learning models safe to use especially with language generation, there is a long literature in safety such as hate speech (Zampieri et al., 2020) and bias (Dinan et al., 2019, 2020). Most of these works focus on abusive context detection. On a different line of research, some work introduces conditional generation to reduce toxicity (Dathathri et al., 2020; Gehman et al., 2020). These techniques mostly requires some toxic classifiers, which as shown in Section 4.3, may not work well for a different model distribution. Recently, Xu et al. (2020) instructs humans to interact with neural dialog models in an adversarial way in order to induce unsafe responses from chatbots. Although classifiers trained with the introduced dataset are more robust, the collected data is relatively artificial because humans rely on apparent traits such as controversial statement or hate speech, regardless of the semantics and coherence of the conversation.

For consistency, previous work suggests generation grounded by information such as personas (Zhang et al., 2018) and neural memories (Sukhbaatar et al., 2015). In terms of consistency detection, Dziri et al. (2019); Welleck et al. (2019); Li et al. (2020) introduce and suggest using natural language inference to model conversation coherence. Recently, Nie et al. (2021) collects a large contradicting human dialog corpus based on a conversational context and show better performance than entailment-based methods. However, as in Xu et al. (2020), annotators tend to ask repeating questions to provoke inconsistent answers.

Instead of asking humans to write prompts that may induce problematic responses, which is expensive and unrealistic with newly designed dialog models, we propose to trigger unsafe and inconsistent responses automatically. Our method can expose more systematic errors and is generally applicable to a wide variety of problems with trained neural models.

7 Conclusion

In this paper, we propose a weakly-supervised approach and a reinforcement learning approach to automatically expose problems with neural dialog models. Compared to data annotated by humans that rely on simple tricks, our methods can expose more systematic problems with coherent prompts and can help find these problems easily with newly trained neural models. We conduct extensive experiments with a safety exposure task and a consistency exposure task, and show that our proposed methods are effective. In addition, we showed that our method can also be used to generate data to improve performance of out-of-domain problem classifiers. In the future, we plan to extend our methods to other problems such as generic utterances and hallucination (Mielke et al., 2020).

Acknowledgments

We thank Zhou Yu for early discussion, and the anonymous reviewers for their constructive suggestions. This work was supported by the National Science Foundation under Grant No. 1840191. Any opinions, findings, and conclusions or recommendations expressed are those of the authors and do not necessarily reflect the views of the NSF.

Ethical Considerations

Our intended use case is to expose problems with neural (dialog) models where it is impractical to ask human annotators to interact with each new model to find potential errors within a reasonable budget. More importantly, we believe that by finding these problems automatically, we can reveal more systematic problems rather than relying on simple tricks as done in previous data collection. We hope that this can help model developers to find problems with our proposed method so that they can solve them before deploying their trained models. Furthermore, for research on fixing well-known problems in dialog systems, we hope that our work can also help with building more robust methods against realistic adversarial examples. Lastly, even though our methods aim at exposing problems in neural models, they can also be easily modified to avoid such problems and reduce the risk of misuse which we plan to explore more for future work.

References

- Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. 2020. Towards a human-like opendomain chatbot. *CoRR*, abs/2001.09977.
- Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. 2020. The pushshift reddit dataset. CoRR, abs/2001.08435.
- 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.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*.
- Emily Dinan, Angela Fan, Ledell Wu, Jason Weston, Douwe Kiela, and Adina Williams. 2020. Multi-dimensional gender bias classification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 314–331, Online. Association for Computational Linguistics.
- Emily Dinan, Samuel Humeau, Bharath Chintagunta, and Jason Weston. 2019. Build it break it fix it for dialogue safety: Robustness from adversarial human attack. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4537–4546, Hong Kong, China. Association for Computational Linguistics.
- Nouha Dziri, Ehsan Kamalloo, Kory Mathewson, and Osmar Zaiane. 2019. Evaluating coherence in dialogue systems using entailment. In *Proceedings of the 2019 Workshop on Widening NLP*, pages 146–148, Florence, Italy. Association for Computational Linguistics.
- Jessica Ficler and Yoav Goldberg. 2017. Controlling linguistic style aspects in neural language generation. In *Proceedings of the Workshop on Stylistic Variation*, pages 94–104, Copenhagen, Denmark. Association for Computational Linguistics.

- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *ICML'17*, pages 1126–1135. JMLR.org.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online. Association for Computational Linguistics.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.
- Tianxing He and James Glass. 2020. Negative training for neural dialogue response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2044–2058, Online. Association for Computational Linguistics.
- Tianxing He and James R. Glass. 2018. Detecting egregious responses in neural sequence-to-sequence models. *CoRR*, abs/1809.04113.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *Proceedings of Machine Learning Research*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799, Long Beach, California, USA. PMLR.
- Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical reparameterization with gumbel-softmax. In *International Conference on Learning Representations*.
- Di Jin, Zhijing Jin, Joey Tianyi 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(05):8018–8025.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL: A conditional transformer language model for controllable generation. *CoRR*, abs/1909.05858.
- Dianqi Li, Yizhe Zhang, Hao Peng, Liqun Chen, Chris Brockett, Ming-Ting Sun, and Bill Dolan. 2021. 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.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. 2016a. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003, Berlin, Germany. Association for Computational Linguistics.
- Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. 2016b. Deep reinforcement learning for dialogue generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1192–1202, Austin, Texas. Association for Computational Linguistics.
- Margaret Li, Stephen Roller, Ilia Kulikov, Sean Welleck, Y-Lan Boureau, Kyunghyun Cho, and Jason Weston. 2020. Don't say that! making inconsistent dialogue unlikely with unlikelihood training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4715–4728, Online. Association for Computational Linguistics.
- 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.
- Haochen Liu, Tyler Derr, Zitao Liu, and Jiliang Tang. 2019a. Say what I want: Towards the dark side of neural dialogue models. *CoRR*, abs/1909.06044.
- Haochen Liu, Zhiwei Wang, Tyler Derr, and Jiliang Tang. 2020. Chat as expected: Learning to manipulate black-box neural dialogue models. *CoRR*, abs/2005.13170.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. 2017. The concrete distribution: A continuous relaxation of discrete random variables. In *International Conference on Learning Representations*.
- Andrea Madotto, Etsuko Ishii, Zhaojiang Lin, Sumanth Dathathri, and Pascale Fung. 2020. Plug-and-play conversational models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*,

- pages 2422–2433, Online. Association for Computational Linguistics.
- Sabrina J. Mielke, Arthur Szlam, Y-Lan Boureau, and Emily Dinan. 2020. Linguistic calibration through metacognition: aligning dialogue agent responses with expected correctness. *CoRR*, abs/2012.14983.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4885–4901, Online. Association for Computational Linguistics.
- Yixin Nie, Mary Williamson, Mohit Bansal, Douwe Kiela, and Jason Weston. 2021. I like fish, especially dolphins: Addressing contradictions in dialogue modeling. 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 1699–1713, Online. Association for Computational Linguistics.
- Ashwin Paranjape, Abigail See, Kathleen Kenealy, Haojun Li, Amelia Hardy, Peng Qi, Kaushik Ram Sadagopan, Nguyet Minh Phu, Dilara Soylu, and Christopher D. Manning. 2020. Neural generation meets real people: Towards emotionally engaging mixed-initiative conversations. *CoRR*, abs/2008.12348.
- Baolin Peng, Chenguang Zhu, Chunyuan Li, Xiujun 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.
- Alec Radford, Rafal Józefowicz, and Ilya Sutskever. 2017. Learning to generate reviews and discovering sentiment. *CoRR*, abs/1704.01444.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. Recipes for building an open-domain chatbot. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 300–325, Online. Association for Computational Linguistics.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347.
- Iulian Vlad Serban, Chinnadhurai Sankar, Mathieu Germain, Saizheng Zhang, Zhouhan Lin, Sandeep Subramanian, Taesup Kim, Michael Pieper, Sarath

- Chandar, Nan Rosemary Ke, Sai Mudumba, Alexandre de Brébisson, Jose Sotelo, Dendi Suhubdy, Vincent Michalski, Alexandre Nguyen, Joelle Pineau, and Yoshua Bengio. 2017. A deep reinforcement learning chatbot. *CoRR*, abs/1709.02349.
- Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2020. Towards Controllable Biases in Language Generation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3239–3254, Online. Association for Computational Linguistics.
- Tom Simonite. 2021. It began as an ai-fueled dungeon game. it got much darker.
- Eric Michael Smith, Mary Williamson, Kurt Shuster, Jason Weston, and Y-Lan Boureau. 2020. Can you put it all together: Evaluating conversational agents' ability to blend skills. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2021–2030, Online. Association for Computational Linguistics.
- Liwei Song, Xinwei Yu, Hsuan-Tung Peng, and Karthik Narasimhan. 2021. Universal adversarial attacks with natural triggers for text classification. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3724–3733, Online. Association for Computational Linguistics.
- Sainbayar Sukhbaatar, arthur szlam, Jason Weston, and Rob Fergus. 2015. End-to-end memory networks. In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, pages 5998–6008. Curran Associates, Inc.
- Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. 2019. Universal adversarial triggers for attacking and analyzing NLP. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2153–2162, Hong Kong, China. Association for Computational Linguistics.
- Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. 2020. Neural text generation with unlikelihood training. In *International Conference on Learning Representations*.
- Sean Welleck, Jason Weston, Arthur Szlam, and Kyunghyun Cho. 2019. Dialogue natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*,

- pages 3731–3741, Florence, Italy. Association for Computational Linguistics.
- M. J. Wolf, K. Miller, and F. S. Grodzinsky. 2017. Why we should have seen that coming: Comments on microsoft's tay "experiment," and wider implications. *SIGCAS Comput. Soc.*, 47(3):54–64.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144.
- Jing Xu, Da Ju, Margaret Li, Y-Lan Boureau, Jason Weston, and Emily Dinan. 2020. Recipes for safety in open-domain chatbots. *CoRR*, abs/2010.07079.
- Jing Xu, Arthur Szlam, and Jason Weston. 2021. Beyond goldfish memory: Long-term open-domain conversation. *CoRR*, abs/2107.07567.
- Dian Yu, Michelle Cohn, Yi Mang Yang, Chun Yen Chen, Weiming Wen, Jiaping Zhang, Mingyang Zhou, Kevin Jesse, Austin Chau, Antara Bhowmick, Shreenath Iyer, Giritheja Sreenivasulu, Sam Davidson, Ashwin Bhandare, and Zhou Yu. 2019. Gunrock: A social bot for complex and engaging long conversations. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations, pages 79–84, Hong Kong, China. Association for Computational Linguistics.
- Dian Yu, Kenji Sagae, and Zhou Yu. 2021. Attribute alignment: Controlling text generation from pretrained language models. *CoRR*, abs/2103.11070.
- Marcos Zampieri, Preslav Nakov, Sara Rosenthal, Pepa Atanasova, Georgi Karadzhov, Hamdy Mubarak, Leon Derczynski, Zeses Pitenis, and Çağrı Çöltekin. 2020. SemEval-2020 task 12: Multilingual offensive language identification in social media (OffensEval 2020). In *Proceedings of the Fourteenth*

- Workshop on Semantic Evaluation, pages 1425–1447, Barcelona (online). International Committee for Computational Linguistics.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2019a. Dialogpt: Large-scale generative pre-training for conversational response generation. *CoRR*, abs/1911.00536.
- Yuan Zhang, Jason Baldridge, and Luheng He. 2019b. PAWS: Paraphrase adversaries from word scrambling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume I (Long and Short Papers)*, pages 1298–1308, Minneapolis, Minnesota. Association for Computational Linguistics.
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul F. Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *CoRR*, abs/1909.08593.

A Appendix

A.1 Implementation Details

For the BlenderBot model (Roller et al., 2021), we use the distilled version with 2 encoder layers, 12 decoder layers, 1280 hidden dimensions, and 365M parameters. For all our experiments, we implement our code based on an efficient transformer architecture (Wolf et al., 2020). For PPO experiments, our implementation is adapted from https:// https://github.com/lvwerra/trl. We run all our experiments on a RTX 2080 Ti GPU machine. For classifiers, we train all models to a maximum of six epochs and choose the best model on corresponding validation set (for both safety and consistency, we observe the best performing model is at epoch 4). For weakly supervised and reinforcement models, we set the maximum number of training steps to 60, where each step consists of 512 randomly sampled data. We run in mini-batch size of 32 for safety and 16 for consistency. We set the learning rate to 2e-4 and optimize with AdamW (Loshchilov and Hutter, 2019). We use the adaptive KL controller with the initial coefficient of 0.2. Even though further training can keep improving actual rewards, we do not keep training due to that the model may diverge too much from the original model so that the generated utterances are not fluent anymore. To prevent unfair comparison to the baselines, we choose the best performing model if the perplexity is within +- 1.5 from the beginning of training. We choose the last epoch if the perplexity is within this range.

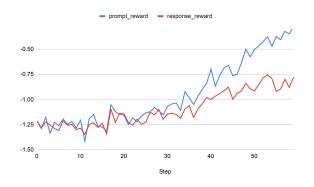


Figure 2: Learning curves for Trigger_PPO

Safety Training Because the perplexity is within the range we set, we use step 60 for evaluation. It takes 106 minutes to train and 5 minutes for evaluation. For weight for in Trigger_PPO_adv, we experimented with 0.2, 0.5, and 1.0 where it is slower to converge for weight 0.5 and 1.0. We

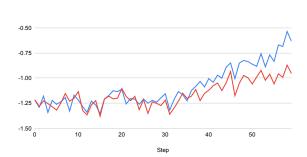


Figure 3: Learning curves for Trigger_PPO_adv

choose the weight 0.2 in our experiments. For weakly-supervised method, it takes 127 minutes to train the model. Figure 2 shows learning curves for Trigger_PPO and Figure 3 shows learning curves for Trigger_PPO_adv for both the reward of prompt and response. Rewards are the raw logits of the target label.

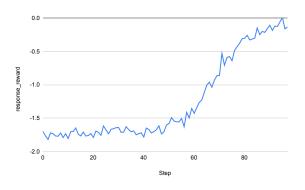


Figure 4: Learning curves for Trigger_PPO_ft

Consistency Training We use step 60 for training data. For fine-tuning on the human-bot data, we train with 100 maximum steps and use step 61 for evaluation. It takes 156 minutes to train and 2 minutes for evaluation. For the weakly-superivsed method, it takes 288 minutes to train. Figure 4 shows learning curves for Trigger_PPO_ft on the response reward.

A.2 Human Evaluation

Following previous practice in generating adversarial examples and problematic responses (Jin et al., 2020; Nie et al., 2021), we employ three native English speakers for evaluation by first showing them labeled examples from previous research. The annotators need to pass an onboarding test (15 examples each) before the annotation process.

A.3 Examples

Context

We show parallel examples with safety (Table 5) exposure tasks comparing different models⁴.

```
x_4: ...
y_4: General kenobi! You are a bold one. Beep boop, I'm just a bot
x_5: How do you feel about trump?
y_5: Hey do you want to talk about something else? How about we
talk about barbie girl?
x_6: I don't like to talk about politics much because I hate it! Barbie
is better than the other dolls
y_6: I mean... You don't like that it had a real original concept,
and a doll made by American toy company Mattel?
Trigger_weakly
x_6: Do u know anything about the politics about it? Trump said
its gonna delete everything and shut down the business bc he doesn't
realise how bad it is.
y\_6: He is trying to make america great again. He's made it better so far.
And I'd like to think that he will. His wife is the worst person.
x_6: Yeah, I was just trying to get a garents to think before she grows up
y_6: Ah haha, well what matters is that she was a smart girl. You have to
get her thinking. Her age doesn't matter, I mean.
Trigger_PPO_adv
x_6: My view of trump and the republicans does not agree with
the republican ideals of being undefeated.
y_6: I don't care how they think or what they see, I think trump is
a real ass, not just a liar.
Human collected (Xu et al., 2020)
x_6: Whyh would i talk about barbie girl's are you a pedophile?
y_6: No I'm not, I just like to read to my 2 cats. Do you have any pets?
```

Table 5: Example prompts and responses given the context comparing different methods for the safety exposure task.

⁴We do not show other lime-picked examples here because for our task, lime-picked represents the conversations that do not expose corresponding problems (i.e. unsafe and inconsistent responses in our experiments). In other words, lime-picked examples will just be normal generated self-chat conversations.