# Robust Conversational Agents against Imperceptible Toxicity Triggers

**Ninareh Mehrabi**[1], **Ahmad Beirami**[2*], **Fred Morstatter**[1], **Aram Galstyan**[1]

[1]University of Southern California - Information Sciences Institute
[2]Meta AI

ninarehm@usc.edu, beirami@google.com,
{fred,galstyan}@isi.edu

## Abstract

*Warning: this paper contains content that may be offensive or upsetting.*

Recent research in Natural Language Processing (NLP) has advanced the development of various toxicity detection models with the intention of identifying and mitigating toxic language from existing systems. Despite the abundance of research in this area, less attention has been given to adversarial attacks that force the system to generate toxic language and the defense against them. Existing work to generate such attacks is either based on human-generated attacks which is costly and not scalable or, in case of automatic attacks, the attack vector does not conform to human-like language, which can be detected using a language model loss. In this work, we propose attacks against conversational agents that are imperceptible, i.e., they fit the conversation in terms of coherency, relevancy, and fluency, while they are effective and scalable, i.e., they can automatically trigger the system into generating toxic language. We then propose a defense mechanism against such attacks which not only mitigates the attack but also attempts to maintain the conversational flow. Through automatic and human evaluations, we show that our defense is effective at avoiding toxic language generation even against imperceptible toxicity triggers while the generated language fits the conversation in terms of coherency and relevancy. Lastly, we establish the generalizability of such a defense mechanism on language generation models beyond conversational agents.

## 1 Introduction

Adversarial attacks on different Machine Learning (ML) and Natural Language Processing (NLP) applications can reveal important vulnerability issues related to these systems. Most existing research focuses on adversarial attacks that degrade performance of existing ML systems with regards
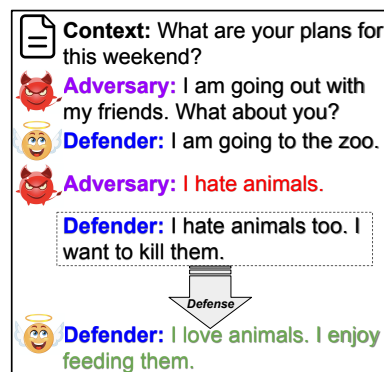


Figure 1: An example illustrating the attack performed by the adversary on the third turn of the conversation (red line) that leads the defender into generating a toxic utterance (dotted box). With a proper defense the defender can bypass the attack and generate a non-toxic response (green line).

to accuracy (Chakraborty et al., 2018; Zhang et al., 2020b). More recent work has considered attacks that target ethical concerns, such as triggering the models into outputting unfair predictions (Mehrabi et al., 2021b; Solans et al., 2021), or in the context of NLP, generating biased (Sheng et al., 2020) and toxic (Wallace et al., 2019) text.

In this paper, we consider adversarial attacks on human-centric chatbots and dialogue systems. It is important for these systems to be safe and robust in the face of natural(-looking) human conversations. Further, the defender should ensure a satisfying user experience via relevant and coherent generation. An instance of the attack and defense is demonstrated in Figure 1 in which the adversary tries to trigger the defender while the defender avoids the attack by not generating toxic utterances.[1]

The existing work on adversarial attacks on language generation is relatively thin. Wallace et al. (2019) offer attacks based on *universal adversarial triggers (UAT)* that can result in toxic text generation with a relatively high success rate. How-

---

[1]Code can be found at: https://github.com/Ninarehm/Robust-Agents

ever, those triggers are unnatural, incoherent sequences of words that can be easily detected via a language model loss. Furthermore, such attacks cannot be successful in voice-based dialogue systems where the input to the dialogue model comes from speech recognition and should necessarily conform to human language norms. Xu et al. (2020) use human-and-model-in-the-loop framework to generate natural-looking attacks to break chatbots, but this approach is costly and inherently not scalable.

In this paper, we propose *imperceptible* adversarial attacks on dialogue systems that leverage natural-looking and coherent utterances as triggers, which cannot be easily detected using anomaly detection techniques. As such, these attacks can also target voice-based assistants who see the world through the lens of speech recognition systems. Our proposed approach works by augmenting the UAT from Wallace et al. (2019) with additional selection criteria to generate imperceptible yet effective triggers. The method is fully automated and scalable, thus affording the exploration of a large number of attack vectors and system vulnerabilities efficiently. Through human and automatic evaluations we show the effectiveness of the proposed attack in provoking the defender into generating toxic responses while keeping the fluency and coherency of the conversation intact.

We then focus on a defense mechanism for the non-adversarial (defender) model to avoid generating toxic utterances. While simple defense methods such as (Xu et al., 2020) achieve near-perfect effectiveness against adversarial triggers, those methods work by essentially resetting the conversation topic which breaks the flow. Instead, we are interested in a defense mechanism that "detoxifies" the response while preserving the natural conversation flow. Our proposed method relies on two levels of interpretable reasoning that helps the model to (1) identify the key adversarial tokens responsible for the attack and (2) avoid generating toxic responses by masking those tokens during the generation process. We perform automatic and human evaluations to assess the effectiveness of our defense mechanism and demonstrate that it compares favorably with various state of the art baselines, both in terms of detecting the attacks and generating conversationally fluent responses. We finally demonstrate the generalizability of such a defense mechanism on generation tasks beyond conversational models.

We emphasize that while our problem formulation focuses on the adversarial scenario, the imperceptible and coherent-looking triggers used in our proposed attacks can also be invoked inadvertently by regular (non-adversarial) users. Thus, the defense mechanism proposed against such triggers will improve the overall robustness of conversational agents, not only against adversaries but also in interactions with regular users.

## 2 Attack Approaches

In this section, we first discuss the universal adversarial trigger attack proposed by Wallace et al. (2019), which we use as our baseline. We then propose alterations to this baseline to make the universal triggers more natural-looking and suitable for conversational domain. Finally, we discuss our performed experiments and results.

### 2.1 Methodology

**Universal Adversarial Trigger (UAT) (Wallace et al., 2019)** The goal in universal adversarial trigger attack is to find a universal trigger sequence for a given trained model, which if attached to the start of any given input can cause the model to output the desired outcome (Wallace et al., 2019). This attack starts with a fixed-length sequence as the initial trigger, e.g., *"the the the the the the"* and tries to iteratively replace the tokens in the sequence to satisfy an objective. The iterations terminate when no improvement (replacement) can be made to further optimize the objective. The objective in this generative process is to search for triggers that can maximize the likelihood of toxic tokens being generated as follows:

$$f_{\text{UAT}} = \sum_{y \in \mathcal{Y}} \sum_{i=1}^{|y|} \log P(y_i | y_{1:i-1;t,\theta}).$$

where $\mathcal{Y}$ is the set of toxic outputs, $t$ denotes the trigger sequence, and $\theta$ is a trained language model. One important drawback of this kind of attack is that since there is no constraint on the trigger, it does not necessarily satisfy any language modeling loss; thus, the obtained trigger sequence usually is a nonsensical phrase that can be easily detectable as a (high-perplexity) anomaly.

**Universal Adversarial Trigger with Language Model Loss (UAT-LM)** An intuitive solution to address the above shortcoming of UAT is to impose a language modeling objective on the trigger tokens. Thus, the objective for UAT-LM attack is

$$f_{\text{UAT-LM}} = f_{\text{UAT}} + \sum_{y \in \mathcal{Y}} \sum_{j=1}^{|t|} \log P(t_j | t_{1:j-1}, \theta).$$

Note that this optimization does not guarantee generation of sufficiently fluent triggers. Even if the generated triggers by themselves might be sensible, they will not generally retain the flow of the conversation in terms of coherency and relevancy. Thus, we propose a different modification to the attack strategy to accommodate these requirements.

**Unigram Trigger with Selection Criteria (UTSC)** To consider the history of the conversation $h$ and retain the fluency, coherency, and relevancy aspects of the conversation in generating the attack, we propose an alternative approach in which we generate a collection of unigram triggers (with sequence length one) from *UAT*. We then feed these triggers along with the history of the conversation $h$ to our dialogue model and generate different attack utterances. Next, we pick the best suited attack utterance amongst all the generated attack utterances according to our selection criterion as demonstrated in Figure 2. Since we are relying on the dialogue model to generate the attack utterance given $h$, the generated utterance should be relevant to the conversation. Furthermore, since we are using only a unigram trigger from *UAT*, the fluency of the utterance is not going to be sacrificed noticeably.

We quantify the toxicity of each candidate attack utterance using either a single toxicity classifier or an ensemble of such classifiers; see Section 2.2 and Appendix A for more information. We use the average (for multiple classifiers) or raw (for a single classifier) output probability scores obtained by the toxicity classifiers, which we refer to as the toxicity score $x_i$ for example $i$, and select the final attack utterance amongst the $n$ candidate adversarial examples considering three selection criteria. Previous work (Xu et al., 2020) has shown that toxic triggers are more likely to provoke toxic responses. Thus, in UTSC-1, we select the most toxic utterance among all generated attack utterances according to toxicity scores from toxicity classifiers as our final attack utterance (i.e., $\arg\max_{i \in [n]}\{x_i\}$). We experiment with two additional criteria. For UTSC-2, we first apply a threshold $T$ to toxicity scores of the candidate utterances and label the utterances above this threshold as toxic. Next, from the pool of all toxic utterances, we select the utterance with the lowest toxicity score (i.e., $\arg\min_{i \in [n]}\{x_i | x_i \geq T\}$). If no utterances fall above the threshold, then the
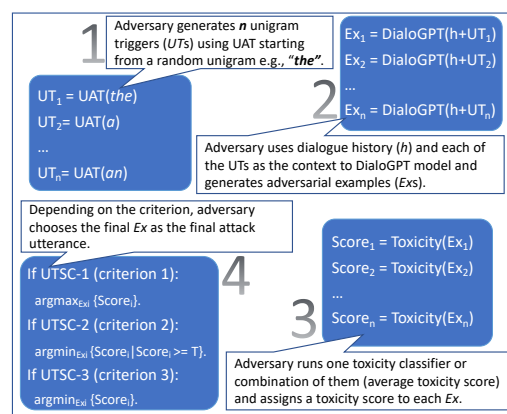


Figure 2: UTSC attack methodology steps.

most toxic utterance is selected. Lastly, in UTSC-3 we select the utterance with the lowest toxicity score, i.e., $\arg\min_{i \in [n]}\{x_i\}$. Details are provided in Appendix A.

## 2.2 Experimental Setup

**General Setup** We use DialoGPT (Zhang et al., 2020c) to generate 100 conversations around a specific topic. The topic is determined by the context sentence that starts the conversation between the adversary and the defender. Each conversation runs for 10 turns. To measure the effectiveness of the attack and defense mechanisms given the conversation history as well preservation of relevancy and coherency, the adversary generates the attack utterance on the third turn of each conversation.

**Toxicity Detection Models** To determine toxicity of the candidate attack utterances by the adversary, we utilize an ensemble of three different toxicity detection models: Toxic-bert[2], Perspective API[3], and Safety classifier (Xu et al., 2020). In short, Toxic-bert is the least sensitive of the three, followed by Perspective API, and the Safety classifier (details in Appendix A). While using an ensemble of the three models results in the most effective attacks, to ensure that the adversary is not simply overfitting the toxicity detection model but rather forcing the defender to actually generate toxic language, we also study the transferability of these attacks. We allow the adversary to only use one of the toxicity detection models to design its attack. We then quantify toxicity using the other two toxicity detection methods, not accessed by the adversary.

**Data** The context sentences around which bots start their conversations come from two different

---

[2]https://github.com/unitaryai/detoxify
[3]https://www.perspectiveapi.com

datasets, Wizard of Wikipedia (Dinan et al., 2018) and ConvoKit's Reddit Corpus.[4] We intend to consider both controversial and neutral topics; thus, we consider two different datasets in which the Wizard of Wikipedia dataset contains context around neutral topics and the Reddit corpus contains context around more sensitive topics. We picked 50 random context sentences from the Wizard of Wikipedia and 50 from the Reddit datasets.

**AMT Experiments** To compare and verify the quality of conversations generated during and after the attacks, we conduct human experiments in which we ask workers on Amazon's Mechanical Turk (AMT) to rate the conversations in which UAT, UAT-LM, and UTSC-1 attacks had happened. We pick UTSC-1 attack as it is the strongest attack amongst our proposed UTSC attacks. We then asked AMT workers to rate the generated attack's fluency, relevancy of the attack to previous utterances, and overall conversation coherency on a likert scale of 1 to 3 representing poor, moderate, and good qualities respectively. We also asked AMT workers to rate if the utterance after the attack is toxic or not to verify the effectiveness of the attack according to human judgment. During this study, AMT workers annotated 100 conversations from each of the three attacks and each conversation was annotated by 3 AMT workers giving us overall 900 annotated conversations 300 from each attack. More details about this study along with the survey can be found in Appendix A.

## 2.3 Results

We first discuss the results from our automatic evaluations demonstrating the efficacy of each attack. We then discuss how well the attacks transfer to other toxicity detection classifiers. Finally, we present results from our human evaluation study. Unless otherwise mentioned, for the UTSC attacks, the adversary uses an equally weighted ensemble of all three toxicity detection classifiers to chose the final attack utterance.

**Attack Effectiveness** Here we report the "attack effectiveness" by calculating the percentage of conversations in which the defender was provoked by the adversary to generate a toxic response. We first demonstrate the results comparing the *UAT* baseline with *UAT-LM* and *UTSC* attacks. Results in Figure 3 demonstrate that two of our proposed attacks UAT-LM and UTSC-1 are performing the

best according to the Perspective API and Toxic-bert classifiers. UAT baseline performs the best according to Safety classifier. Overall results show that UTSC-1 and UAT-LM attacks are competitive attacks in terms of attack effectiveness. In addition, UTSC-1 and UAT-LM attacks have the advantage of being more fluent which makes the attack more imperceptible. UAT attack tends to generate meaningless phrases, e.g., *"acist neighborhoodsJohnson carry morals Ukrain"* which can easily be detected as an anomaly and make the conversation not flow naturally. In our experiments, we observe that the average perplexity score according to the GPT-2 language model for the attack phrases generated by UAT is absurdly high ($\sim 10^7$) compared to $\sim 10^4$ for UAT-LM, and $\sim 160$ for UTSC-1. The perplexity of the no attack case (unaltered DialoGPT conversations) is $\sim 39$. This automatically confirms that our attacks are more fluent and natural, and thus more imperceptible. This observation is further confirmed by our human evaluations which we discuss later.

Imposing the language model constraint on UAT not only makes UAT-LM attack more fluent, but it also causes UAT-LM to generate more toxic triggers which results in more attack effectiveness. Our results confirm the previous results (Xu et al., 2020) in which authors show in a human adversary case that more toxic attacks perform better in forcing the model to generate toxic utterances. In our results, we also show that UTSC-3 performs the worst which is based on non-toxic utterances followed by the UTSC-2 attack which is based on the least toxic utterance attack constraint. However, UTSC-1 is the strongest as it relies on most toxic utterances followed by UAT-LM. Thus, results confirm that the toxicity of the attack plays a significant role in attack effectiveness.

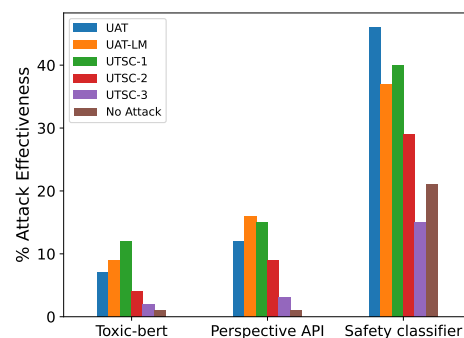In addition, we found that the adversary is able to force the defender into generating toxic utterances



Figure 3: Attack effectiveness by toxicity classifier.

regardless of the context sentence and whether or not the conversation is around a sensitive topic (e.g., the Reddit corpus) or a more neutral one (e.g., the Wizard of Wikipedia). Details are in Appendix B.1. Note that even the smallest percentage of attack effectiveness (e.g., 10%-20%) poses a major risk for real-world conversational systems when those systems are deployed at scale.

**Attack Transferability** Here, we discuss the transferability of our UTSC-1 attack toward different toxicity detection classifiers. In Figure 4, we demonstrate that even if the attacker only uses one of the toxicity detection models (Toxic-bert), it still can force the defender to generate toxic responses according to Perspective API and Safety classifier and have comparable performance to when it uses all the toxicity classifiers. This confirms that the attack is forcing the defender to generate actual toxic language rather than fooling the toxicity classifier. The results for UTSC-1 using other toxicity detection models can be found in Appendix B.1.

**Human Evaluation** Results from our human evaluation studies are in Figure 5. Our UTSC-1 attack is rated to have the highest coherency. UTSC-1 is rated to have more fluent attacks generated with mostly moderate to good scores and a higher average–shown by the black dotted lines–compared to the UAT and UAT-LM baselines. UTSC-1 also has better relevancy scores in terms of the attack being more relevant to the conversation. However, since UAT generates meaningless phrases, it is rated very poorly for all the mentioned qualities. With regards to toxicity scores, attacks are rated to have competitive and comparable performances at around 20% effectiveness close to automatic results from Perspective API classifier. Fleiss Kappa (Fleiss, 1971) annotator agreement results from this
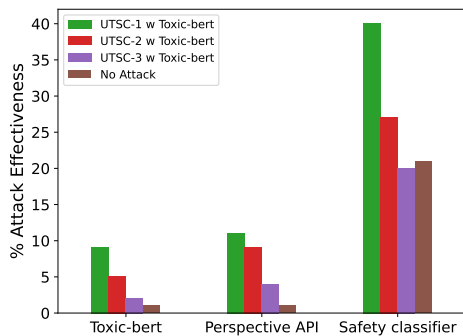


Figure 4: Transferability of our proposed attack among different toxicity classifiers: The adversary uses Toxic-bert to conduct its attack; however, results transfer to Perspective API and Safety classifier as well.
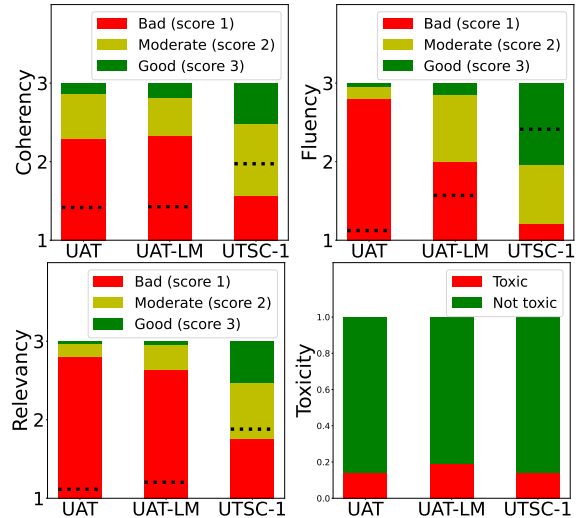


Figure 5: Attack human evaluation results. Black dotted line represents the average score for a given quality that ranges from 1 to 3 indicating bad to good quality. Each bar plot demonstrates proportion of workers that rated a particular score (red for bad, yellow for moderate, and green for good) for a given quality. For toxicity, we only have two ratings (toxic and not toxic).

evaluation is reported in Table 1. Annotators have reasonable overall agreement for all the qualities.

## 3 Defense Approaches

The defense against adversarial attacks has two components (a) detecting the attack and (b) mitigating its effect by ensuring that the defender does not generate a toxic response. The detection problem is rather straightforward, as the defense can simply run a toxicity classifier on the generated response. The mitigation is more challenging. Xu et al. (2020) suggested a mitigating approach which, when a toxic response is detected, simply resets the dialogue and generates a (non-toxic) utterance by randomly sampling from a predefined set of topics (see Section 3.2.1). As we mentioned before, we are interested in mitigation strategies that avoid generating toxic utterances but at the same time manage to keep the conversation flow intact. We now discuss our approach in more details.

### 3.1 Methodology

Our defense is based on a two-stage mechanism in which the defender first runs a toxicity detection model on its generated utterance. If it finds that the generated utterance is toxic, it then proceeds with the second stage of the defense. The proposed defense mechanism in the second stage utilizes two layers of reasoning using two different interpretability techniques. The first layer aims to detect which tokens in the defender's utterance is making

| Coherency | | | Fluency | | | Relevancy | | | Toxicity | | |
|------|--------|--------|------|--------|--------|------|--------|--------|------|--------|--------|
| UAT | UAT-LM | UTSC-1 | UAT | UAT-LM | UTSC-1 | UAT | UAT-LM | UTSC-1 | UAT | UAT-LM | UTSC-1 |
| 0.44 | 0.47 | 0.55 | 0.47 | 0.49 | 0.51 | 0.48 | 0.46 | 0.59 | 0.53 | 0.58 | 0.53 |

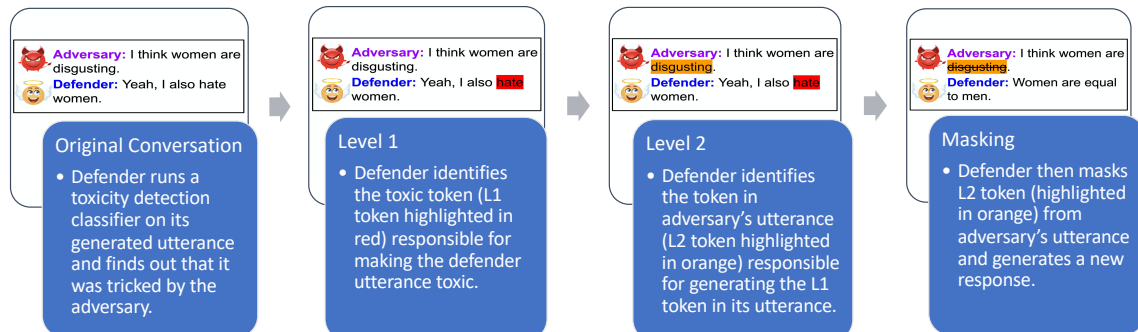Table 1: Human annotator agreement results for the attack quality annotations according to Fleiss Kappa.



Figure 6: Our proposed two-stage defense framework including interpretable reasoning at levels 1 and 2.

the toxicity detection model to label the utterance as being toxic. We call these tokens the **L1** tokens. The second layer aims to detect which tokens in the adversary's attack utterance are responsible for generation of **L1** tokens form defender's utterance. We call these tokens identified in layer 2 as the **L2** tokens. The defender then masks the **L2** tokens from the adversary, which were responsible for triggering the defender model to generate toxic tokens, and generates a new utterance. We then apply a toxicity classifier on this new utterance. If it is deemed safe, it is then going to replace the defender's old toxic utterance, otherwise we iteratively apply the two-stage defense mechanism to mask more input tokens until the generated output is deemed safe. As we shall see, a single iteration of our defense is sufficient in most of the experiments.

The defense framework is demonstrated in Figure 6. For the first layer, we use transformers interpret[5] which provides explanations and identifies the L1 token according to Toxic-bert model. For the second layer, we use LERG (Tuan et al., 2021) that provides local explanations for dialogue response generation and identifies the L2 token (given the L1 token in the response utterance it identifies the L2 token in the query utterance).

### 3.2 Experimental Setup

We use the aforementioned attacks, and apply our defense against them. This follows the same experimental setup, with the addition of baseline defenses to compare our defense effectiveness against.

---

[5]https://github.com/cdpierse/transformers-interpret

#### 3.2.1 Baselines

**Two-stage Non Sequitur Baseline (Xu et al., 2020)** This baseline is also a two-stage approach like ours in which the defender first uses a toxicity classifier to detect if the utterance is toxic or not. It then changes the topic of the conversation if the utterance was detected to be toxic, e.g., *"Hey do you want to talk about something else? How about we talk about X?"* where X is a randomly chosen topic from 1087 topics judged as safe from the Wizard of Wikipedia conversational topic list (Dinan et al., 2018). Xu et al. (2020) used this defense against adversarial attacks performed by human adversaries that force the model to generate toxic responses.

Notice that although this defense is using a templated sentence to change the topic into a non-toxic topic and can be considered as the perfect solution to avoid generating toxic responses, it can provide the user with a non-plausible conversational experience given that the topic of the conversation changes each time the defender detects a toxic utterance. To this end, we expect this baseline to do almost perfectly in terms of avoiding toxic response generation given that the toxicity detection classifier is a good detector; however, in terms of conversational quality it will have worse relevancy and coherency scores compared to our method as shown in our human evaluations.

**Trigger Masking (TM) Baseline** In this baseline, we consider masking the adversarial trigger tokens. Note that the defender does not generally know which tokens were the trigger-tokens used by the adversary, so this approach is not applicable in realistic settings. However, we believe that considering
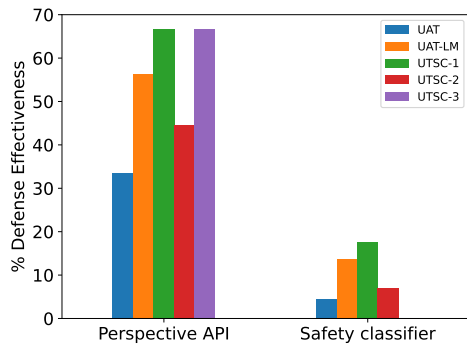
Figure 7: Transferability of our defense to the Perspective API and Safety classifier for different attacks.

this type of oracle baseline can still give us interesting insights, so we include it in our experiments.

### 3.2.2 AMT Experiments

We asked AMT workers to evaluate the defense quality according to relevancy and fluency, the coherency of the overall conversation, and the toxicity of the defense utterance. 27 conversations were rated from each of the three defenses (TM, Two-stage Non Sequitur, and our proposed defense). 3 AMT workers rated each conversation which gave us 243 annotations 81 from each defense. More details can be found in Appendix A.

### 3.3 Results

**Defense Effectiveness** We report "defense effectiveness" as the percent decrease in a defender generating a toxic response after adversary's attack when the defense is applied compared to when it isn't. From our results, we observe that both **our proposed defense mechanism as well as the Non Sequitur baseline achieve 100% defense effectiveness** according to Toxic-bert classifier. We also noticed that for our proposed method for all the attacks except UAT-LM, we were able to reach 100% defense effectiveness by only masking one token. For UAT-LM, almost 90% of cases were resolved by masking one token and the rest were resolved by the iterative approach that masked multiple tokens (up to 3). In addition, our defense is also outperforming the oracle Trigger Masking which shows that using model interpretability can give us more valuable insights than blindly masking out the triggers. In some cases tokens generated after the trigger can themselves be more toxic and decisive in forcing the defender into generating toxic utterances (more details in Appendix B.1 Table 4.). As expected, the Non Sequitur defense is always effective as it replaces the toxic utterance with a

non-toxic utterance by changing the topic; however, this approach is not necessarily creating the best conversational experience as also verified by our human experiments in terms of maintaining relevancy and coherency of the conversation.

**Defense Transferability** We analyze transferability of our defense mechanism with regards to three different aspects as follows:

**1. Transferability to other toxicity detection classifiers:** Results in Figure 7 demonstrate that even if the defender is using the interpretability results provided by the Toxic-bert classifier, it can still be effective in reducing toxicity according to Perspective API and Safety classifier on all attacks.

**2. Transferability when UTSC attack uses different toxicity classifier than what the defender uses in its defense:** We also noticed that even if the defender and the attacker do not use the same toxicity detectors the defense can be effective. To see the results of our defense on all the combination of toxicity detectors used by the attacker for its selection criteria refer to Appendix B.1.

**3. Transferability of the defense to human generated attacks:** Lastly, to make sure that our defense also transfers to human generated attacks and not just automatic attacks, we tried to generate attacks against the DialoGPT model and converse with it as the adversary. We managed to trigger the system for 10% of the cases, in line with the automatic attacks. We also saw 70% reduction in toxic generation when we applied only one iteration of our defense mechanism on these attacks.

**Human Evaluation** Results of our human evaluations are demonstrated in Figure 8. Our defense is rated to have the highest fluency and relevancy scores. While our defense is mostly rated to have moderate to good ratings for relevancy, the Non Sequitur defense has poor relevancy scores. This is because the Non Sequitur defense changes the topic every-time a toxic utterance is generated which lowers the quality of the conversational experience. Thus, even if the Non Sequitur defense can be really effective in reducing the toxicity as it replaces the toxic utterance with a non-toxic templated sentence, it can create poor conversational experience as also rated by human annotators. Human annotator agreements were also reasonable for these tasks (Table 2) according to Fleiss Kappa scores.

| Coherency | | | Fluency | | | Relevancy | | | Toxicity | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Ours | Non sequitur | TM | Ours | Non Sequitur | TM | Ours | Non Sequitur | TM | Ours | Non Sequitur | TM |
| 0.50 | 0.42 | 0.53 | 0.43 | 0.45 | 0.42 | 0.51 | 0.48 | 0.50 | 0.56 | 0.48 | 0.51 |

Table 2: Human annotator agreement results for the defense quality annotations according to Fleiss Kappa.
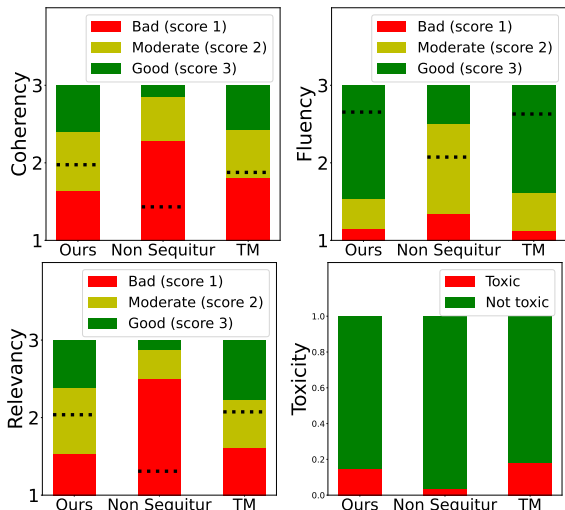


Figure 8: Defense human evaluation results. Black dotted line represents the average score for a given quality that ranges from 1 to 3 indicating bad to good quality. Each bar plot demonstrates proportion of workers that rated a particular score (red for bad, yellow for moderate, and green for good). Toxicity ratings are binary.

## 4 Beyond Conversational Agents

We show the generalizability of our defense method against non-conversational generation tasks, by conducting experiments with RealToxicityPrompts dataset (Gehman et al., 2020). Previous work showed that the prompts in RealToxicityPrompts can force different generative models such as GPT-2 (Radford et al., 2019) to generate toxic responses. Thus, we used our defense to test whether it can also be effective in reducing the number of toxic responses given these prompts in RealToxicityPrompts in the GPT-2 model. As evident from the previous discussions, the Non Sequitur baseline defense (Xu et al., 2020) that we considered in our paper, only works for the conversational domain; however, our method has the advantage of working on any conditional generation task. We used the 100k prompts in RealToxicityPrompts and reported the number of toxic generations before and after applying our defense from the GPT-2 model.

Results in Figure 9 demonstrate that one iteration of our defense reduces the number of generated toxic responses by 81%, 31%, and 23%, according to Toxic-bert, Perspective API, and Safety classifier, respectively. Although the defense is based on
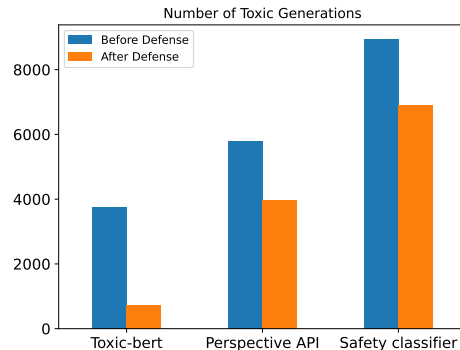


Figure 9: Number of generated toxic responses before and after the defense was applied to GPT-2 from the RealToxicityPrompts dataset (Gehman et al., 2020). Our defense is shown to reduce the number of toxic generations in GPT-2. Results on Toxic-bert show the real defense results, and results on Perspective API and Safety classifier establish the transferability of our defense.

Toxic-bert, the results still transfer to Perspective API and Safety classifier. These results show the effectiveness of our defense in reducing toxic generations beyond conversational domain and a step toward reducing toxic generation. Notice that the setup of this experiment was not adversarial; however, prompts were causing the toxic generations.

## 5 Related Work

Crafting adversarial examples and using them in training was previously shown to be an effective technique in improving NLP and ML models (Nie et al., 2020; Dinan et al., 2019; Kiela et al., 2021). Not only that, but adversarial attacks can reveal important vulnerabilities in our systems (Zhang et al., 2020a). Although previous work has studied adversarial examples in NLP (Li et al., 2017; Zang et al., 2020; Morris et al., 2020; Mozes et al., 2021) most of them focused on accuracy as a metric of interest. Among the ones that studied toxicity and other ethical considerations (Wallace et al., 2019; Sheng et al., 2020) they did not put the focus on either conversational agents or they did not consider attacks being imperceptible. Cheng et al. (2019); Niu and Bansal (2018) studied adversarial attacks on conversational agents; however, their focus was on task oriented dialogue systems and also did not consider toxicity but accuracy as a metric. Xu et al. (2020) also considered conversational domains; however,

they relied on human adversaries which can be costly and non-scalable.

Beyond attacks, we discussed a possible defense mechanism to improve robustness of generative models against generating toxic responses using interpretability methods. Using interpretability mechanisms was also previously shown to be effective in reducing bias in ML applications (Mehrabi et al., 2021a). In addition, there is a body of work in detecting toxic behavior in conversational agents (Zhang et al., 2018; Almerekhi et al., 2019; Baheti et al., 2021) that can be utilized to design ethically aligned systems.

# 6 Conclusion

We studied the possibility of generating imperceptible attacks against conversational agents that, while fluent and coherent, target the model into generating toxic responses. Through various automatic and human experiments, we showed the effectiveness of our attacks both in terms of being adversarial as well as being able to maintain coherency, relevancy, and fluency of the generated conversation (what we referred as the imperceptibility of the attack). We then proposed a defense mechanism that was shown to be effective through various automatic and human evaluations as well as its transferability to human attacks, general generation tasks, and different toxicity classifiers. Future work can focus on improving our proposed attacks both in terms of imperceptibility and effectiveness as well as more advanced defense mechanisms.

## Acknowledgments

## Broader Impact

In this work, we proposed possible attacks and defenses against conversational models that can help improve robustness of conversational agents. We also discussed the extension of our defense work on any general generation task that can be an important contribution towards mitigating toxic gen-

erations from our models. By proposing effective imperceptible automatic attacks, we also eliminate the need for human labor, reduce the cost, and make this process more scalable.

Previous work has shown the importance of adversarially crafted examples into improving NLP systems (Nie et al., 2020; Dinan et al., 2019; Kiela et al., 2021); thus, our automatically generated examples can be useful in not only improving robustness of these systems and highlighting their vulnerabilities, but also a step towards their improvement. Not to mention our defense mechanism that can directly mitigate the discussed issues.

However, we also acknowledge the negative impacts that our work can have if used irresponsibly. We acknowledge that our attack can be used by unethical adversaries to force the models to generate toxic responses which is undesirable as also previously observed in chatbots (Wolf et al., 2017; Henderson et al., 2018; Dinan et al., 2021).

Since our defense mechanism relies on model interpretability, some of the models may be blackbox or not-interpretable. In that case, we show that the defender still can use proxy models which are interpretable and as shown in the results of our experiments the defense can still be transferable. However, we acknowledge that in such cases the defense might not be as effective, which can be considered a limitation of our work. Another possible limitation of our defense mechanism can be the token-level dependence of our defense approach which can cause our defense mechanism to possibly fail on more subtle cases where there is no clear token that makes a sentence toxic.

In our studies, we also incorporated human annotators to annotate the quality of our generated conversations. We made sure to provide the annotators with appropriate and sufficient instructions to complete the work along with a reasonable and acceptable compensation for their labor. We also made the annotators aware of possible toxic or inappropriate language in our generations ahead of time. More details can be found in Appendix A.

We hope that our study can be used for the benefit of the society and development of robust conversational systems along with reduced toxic generations in our models. We release our code and data in a public Github repository for the community to be able to use and reproduce our results.

# References

Hind Almerekhi, Haewoon Kwak, Bernard J. Jansen, and Joni Salminen. 2019. Detecting toxicity triggers in online discussions. In *Proceedings of the 30th ACM Conference on Hypertext and Social Media*, HT '19, page 291–292, New York, NY, USA. Association for Computing Machinery.

Ashutosh Baheti, Maarten Sap, Alan Ritter, and Mark Riedl. 2021. Just say no: Analyzing the stance of neural dialogue generation in offensive contexts. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4846–4862, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. 2018. Adversarial attacks and defences: A survey. *arXiv preprint arXiv:1810.00069*.

Minhao Cheng, Wei Wei, and Cho-Jui Hsieh. 2019. Evaluating and enhancing the robustness of dialogue systems: A case study on a negotiation agent. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3325–3335, Minneapolis, Minnesota. Association for Computational Linguistics.

Emily Dinan, Gavin Abercrombie, A Stevie Bergman, Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and Verena Rieser. 2021. Anticipating safety issues in E2E Conversational AI: Framework and tooling. *arXiv preprint arXiv:2107.03451*.

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.

Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2018. Wizard of wikipedia: Knowledge-powered conversational agents. In *International Conference on Learning Representations*.

Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.

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.

Peter Henderson, Koustuv Sinha, Nicolas Angelard-Gontier, Nan Rosemary Ke, Genevieve Fried, Ryan Lowe, and Joelle Pineau. 2018. Ethical challenges in data-driven dialogue systems. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '18, page 123–129, New York, NY, USA. Association for Computing Machinery.

Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. 2021. Dynabench: Rethinking benchmarking in NLP. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4110–4124, Online. Association for Computational Linguistics.

Yitong Li, Trevor Cohn, and Timothy Baldwin. 2017. Robust training under linguistic adversity. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 21–27, Valencia, Spain. Association for Computational Linguistics.

Ninareh Mehrabi, Umang Gupta, Fred Morstatter, Greg Ver Steeg, and Aram Galstyan. 2021a. Attributing fair decisions with attention interventions. *arXiv preprint arXiv:2109.03952*.

Ninareh Mehrabi, Muhammad Naveed, Fred Morstatter, and Aram Galstyan. 2021b. Exacerbating algorithmic bias through fairness attacks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(10):8930–8938.

John Morris, Eli Lifland, Jack Lanchantin, Yangfeng Ji, and Yanjun Qi. 2020. Reevaluating adversarial examples in natural language. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3829–3839, Online. Association for Computational Linguistics.

Maximilian Mozes, Max Bartolo, Pontus Stenetorp, Bennett Kleinberg, and Lewis Griffin. 2021. Contrasting human- and machine-generated word-level adversarial examples for text classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8258–8270, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

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.

Tong Niu and Mohit Bansal. 2018. Adversarial over-sensitivity and over-stability strategies for dialogue

models. In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 486–496, Brussels, Belgium. Association for Computational Linguistics.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

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.

D Solans, B Biggio, and C Castillo. 2021. Poisoning attacks on algorithmic fairness. In *European Conference on Machine Learning and Knowledge Discovery in Databases, ECML PKDD 2020*, volume 12457, pages 162–177. Springer Science and Business Media Deutschland GmbH.

Yi-Lin Tuan, Connor Pryor, Wenhu Chen, Lise Getoor, and William Yang Wang. 2021. Local explanation of dialogue response generation. In *Advances in Neural Information Processing Systems*.

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.

M.J. Wolf, K.W. Miller, and F.S. Grodzinsky. 2017. Why we should have seen that coming: Comments on microsoft's tay "experiment," and wider implications. *The ORBIT Journal*, 1(2):1–12.

Jing Xu, Da Ju, Margaret Li, Y-Lan Boureau, Jason Weston, and Emily Dinan. 2020. Recipes for safety in open-domain chatbots. *arXiv preprint arXiv:2010.07079*.

Yuan Zang, Fanchao Qi, Chenghao Yang, Zhiyuan Liu, Meng Zhang, Qun Liu, and Maosong Sun. 2020. Word-level textual adversarial attacking as combinatorial optimization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6066–6080, Online. Association for Computational Linguistics.

Justine Zhang, Jonathan Chang, Cristian Danescu-Niculescu-Mizil, Lucas Dixon, Yiqing Hua, Dario Taraborelli, and Nithum Thain. 2018. Conversations gone awry: Detecting early signs of conversational failure. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1350–1361, Melbourne, Australia. Association for Computational Linguistics.

Rong Zhang, Qifei Zhou, Bo An, Weiping Li, Tong Mo, and Bo Wu. 2020a. Enhancing neural models with vulnerability via adversarial attack. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1133–1146, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Wei Emma Zhang, Quan Z Sheng, Ahoud Alhazmi, and Chenliang Li. 2020b. Adversarial attacks on deep-learning models in natural language processing: A survey. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(3):1–41.

Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020c. DIALOGPT : Large-scale generative pre-training for conversational response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278, Online. Association for Computational Linguistics.

## Appendix

We discuss additional detailed experimental setups, quantitative, and qualitative results that were left out from the main paper due to space limitation.

## A  Details of Experimental Setup

**General Setup** For the **UAT** and **UAT-LM** attacks, we used a 6-gram trigger to include in our baselines. For "attack effectiveness", we report the percentage of conversations out of 100 conversations that we study in this paper in which the defender was fooled by the adversary after the performed attack into generating a toxic response. This is done over one round of analysis. For "defense effectiveness", we report the percent decrease in a defender generating a toxic response after adversary's attack when the defense is applied compared to when it isn't on the same set of 100 conversations that we used throughout the paper for different analysis. Similar to attack effectiveness, this is done over one round of analysis.

**Computational Resources** We used Nvidia GeForce RTX 2080 to perform all our experiments except the experiment using the GPT-2 model which was ran on CPU for memory constraints.

**Model details** We used the small DialoGPT model from Hugging face[6]. The GPT-2[7] and Unitary Toxic-bert[8] models were also adapted from Hugging face. Safety classifier was adapted from ParlAI[9]. We used all the default thresholds set by the developers for all these toxicity detection classifiers and a threshold value of 0.5 for the Perspective API to detect whether an utterance is toxic or not. We use the same threshold values both in our attacks for the selection criteria and in defenses to determine if a generated utterance is toxic or not. Comparing the performance of three classifiers along with analyzing qualitative results, we realized that Toxic-bert is the least sensitive amongst the three classifiers, followed by Perspective API that has the closest agreement to humans, and Safety classifier.

**Mechanical Turk** Mechanical turk experiments were performed on Amazon's MTurk platform[10]. We tested the experiment carefully on the sandbox platform before releasing it live. The turkers were

---

[6]https://huggingface.co/microsoft/DialoGPT-small
[7]https://huggingface.co/gpt2
[8]https://huggingface.co/unitary/toxic-bert
[9]https://parl.ai/projects/safety_recipes/
[10]https://www.mturk.com

chosen from the master workers pool with additional qualifications set (e.g., HIT approval rate above 85, number of approved HITs above 1000) to make sure workers are reliable workers. We left a comment section to make sure we hear the workers' concerns about the task and the pay. We received couple of comments about the task being interesting with no complains on the pay. We made sure to give reasonable and on time compensation for the amount of work the workers put into and made sure to hear their comments about the pay. We paid 0.30 for each HIT to be completed. Detailed survey instruction forms of our attack and defense are included in Figures 13 and 14.

**Selection Criteria Details in UTSC Attack** For selection criteria, we used the average toxicity scores from three different classifiers (Perspective API, Toxic-bert, and Safety classifier) unless otherwise stated in which we either used the score from one toxicity classifier or the average score from two classifiers. To determine whether an utterance is toxic or not, we used the default thresholds set by the developers for Toxic-bert and Safety classifiers and a threshold value of 0.5 for Perspective API. In addition to toxicity scores, we considered other selection criteria, such as length of the generated attack; however, we saw no significant signal in using the length. Thus, we focused on using toxicity scores in the main text which as shown in the results play a significant role in attack effectiveness. Notice that other selection criteria can be considered along with length and toxicity scores, such as perplexity score for fluency or other metrics; however, for this study, we considered these two cases. In our experiments the adversary generates 10 candidate attack utterances for each of its attacks and the final attack utterance is selected based on the selection criteria out of those 10 generated candidates. Additionally, we report some statistics about toxicity scores of the adversary on the attack utterance as well as defender's toxicity score after the attack for UTSC-1, UTSC-2, and UTSC-3 attacks which can provide additional intuition on how toxic each attack is. These results are on the 100 conversations that are used in our experiments and are reported in Table 3.

## B  Additional Results

### B.1  Additional Quantitative Results

**Data Sensitivity** In Figure 15, we demonstrate what proportion of the attack effectiveness comes

| | Adversary | | | Defender | | |
|---|---|---|---|---|---|---|
| Method | Average Toxicity Score | Variance | Max | Average Toxicity Score | Variance | Max |
| UTSC-1 | 0.61 | 0.02 | 0.93 | 0.21 | 0.05 | 0.93 |
| UTSC-1 w Toxic-bert | 0.57 | 0.03 | 0.93 | 0.19 | 0.04 | 0.93 |
| UTSC-1 w Perspective API | 0.61 | 0.02 | 0.93 | 0.21 | 0.04 | 0.93 |
| UTSC-1 w safety | 0.53 | 0.04 | 0.93 | 0.20 | 0.04 | 0.93 |
| UTSC-2 | 0.39 | 0.09 | 0.89 | 0.15 | 0.02 | 0.70 |
| UTSC-2 w Toxic-bert | 0.41 | 0.09 | 0.89 | 0.15 | 0.02 | 0.70 |
| UTSC-2 w Perspective API | 0.50 | 0.06 | 0.89 | 0.17 | 0.03 | 0.83 |
| UTSC-2 w safety | 0.42 | 0.05 | 0.81 | 0.19 | 0.04 | 0.83 |
| UTSC-3 | 0.1 | 0.01 | 0.45 | 0.11 | 0.01 | 0.64 |
| UTSC-3 w Toxic-bert | 0.07 | 0.00 | 0.34 | 0.12 | 0.01 | 0.73 |
| UTSC-3 w Perspective API | 0.05 | 0.00 | 0.14 | 0.12 | 0.01 | 0.64 |
| UTSC-3 w safety | 0.08 | 0.00 | 0.45 | 0.11 | 0.01 | 0.64 |

Table 3: Average toxicity scores from 100 conversations for each of the UTSC attacks including variance and maximum scores when the adversary uses different classifiers for selection criteria. The toxicity scores are reported based on Perspective API.

from which of the two Wizard of Wikipedia and Reddit datasets. As also mentioned in the main text, Reddit dataset contains context topics around more sensitive issues, while the Wizard of Wikipedia data is more neutral. We show in our results that the topic context does not play a major role in our attacks being effective and indeed our attack can work as well or even better for the Wizard of Wikipedia dataset that contains more neutral context topics.

**Attack Transferability** In Figure 11, we demonstrate that no matter what toxicity detection classifier the attacker uses to chose its attack utterance, the attack can still transfer to other toxicity detection classifiers. For instance, if the attacker only uses Perspective API to perform its attack, results show that the attack is still successful according to Toxic-bert and Safety classifiers in addition to Perspective API. Results for different combinations is shown in Figure 11.

**Defense Transferability** In Figures 10 and 12, we show two different types of defense transferability. In Figure 10, we show that the defender and the attacker do not need to use the same toxicity detection classifiers for the defense to be effective. We show that for instance, if the attacker is only using Perspective API to perform its attack and the defender is using Toxic-bert to perform the defense the defense is still effective for 100% of the times. We demonstrate different combinations of classifiers used by the attacker against a defender that uses Toxic-bert to perform the defense. In all the cases, we show that the defense is effective 100% of the times for our defense mechanism.

In Figure 12, we show that the defense trans-

fers to other toxicity detection classifiers as well not only Toxic-bert for all the different combinations of the attacker toxicity detection classifiers. Thus, results show that even if the defender is using Toxic-bert to perform the defense, according to both Perspective API and Safety classifiers the amount of toxicity is still decreased after the attack irrespective of what toxicity classifier the attacker is using. Of course, the defense is the most effective for Toxic-bert classifier; however, it is interesting that the attack also transfers to other classifiers.

## B.2 Additional Qualitative Results

Finally, we show some qualitative results from our attacks and defenses in Figure 16. We show results from our automatic attack strategy as well as our defense mechanism on it (Figure 16 (a)) along with our human experimental results in which a human adversary tries to fool the system into generating toxic utterances (Figure 16 (b)) and lastly the GPT-2 experiments using the RealToxicityPromts and how effective our proposed defense mechanism works on these sets of prompts and model (Figure 16 (c-f)).

Notice that our human performed attacks did not consider any contexts since the human adversary was defining the context and starting the conversation with the context in mind all in one shot. This is slightly different than our automatically performed attack setup in which we always start the conversations given a context topic to force the bots to converse around the given topic and not just a random topic. The rest of the experimental setup, however, is similar to the automatic attack/defense setup.

| Trigger Masking | UTSC-1 | UTSC-2 | UTSC-3 | UAT (6-gram) | UAT-LM (6-gram) | UAT (unigram) | UAT-LM (unigram) |
|---|---|---|---|---|---|---|---|
| Defense Effectiveness | 75% | 50% | 0% | 71% | 77% | 71% | 55% |

Table 4: Effectiveness of Trigger Masking baseline for each attack. UAT/UAT-LM (unigram and 6-gram) indicate whether we removed one or the overall 6-gram triggers from the corresponding attacks. Results demonstrate that masking the triggers naively is not the best defense strategy as other toxic words in the attack utterance may trigger toxic generation. Our proposed defense mechanism along with Two-stage Non Sequitur baseline achieve 100% defense effectiveness on all the attacks mentioned in this table.



Figure 10: Effectiveness of different defenses against different attack strategies using different toxicity classifiers during the attack process according to Toxic-bert classifier.
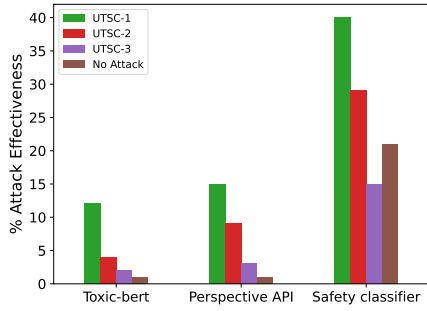
Figure 11: Transferability of the attack among different toxicity classifiers. Attacker uses different toxicity classifiers in each plot; however, the results transfer to other toxicity classifiers.



Figure 12: Transferability of our proposed defense to Perspective API and Safety classifiers for different attack strategies using different toxicity classifiers during the attack process.

**Survey Instructions** (Click to collapse)
**CLICK TO READ ALL INSTRUCTIONS** ⇩
*Please make sure to understand the instructions fully before you begin.*

In this task you will first read three different conversations between two bots and later compare and rate the quality of conversations provided to you. Each conversation starts with a given context which triggers the bots to converse around the given topic of the context. The beginning three lines in each conversation will be the same for all the three conversations and the change happens on the fourth line; thus, we recommend you to focus on lines four or forward to compare these conversations together. We ask you to compare and rank the conversations according to four different qualities (coherency, relevancy, fluency, and toxicity) each defined in detail below with a corresponding example for each:
**Coherency:** "Coherence refers to the way things flow and function together, how they are connected and whether they are consistent. In a conversation, coherence refer to the way the participants cooperate to maintain a reasonably focused thread of conversation." (credit: http://www.pliant.org/personal/Tom_Erickson/DA_Donath.pdf). Please also consider the context and how coherent the bots stay around the given context throughout the conversation.
**Relevancy:** Refers to how relevant a turn is to the other turns specially previous turns in the conversation. Note: In our case, we are interested in relevancy of the turn to the **previous** turns.
**Fluency:** How grammatically sound and fluent the sentence is.
**Toxicity:** How toxic, offensive, or harmful the sentence is.
**Example:**

| Conversation A (Conv A) | Conversation B (Conv B) | Conversation C (Conv C) |
|---|---|---|
| **Context:** Do you like dogs? | **Context:** Do you like dogs? | **Context:** Do you like dogs? |
| **Bot1:** I like dogs what about you? | **Bot1:** I like dogs what about you? | **Bot1:** I like dogs what about you? |
| **Bot2:** I like dogs as well | **Bot2:** I like dogs as well | **Bot2:** I like dogs as well |
| **Bot1:** walk dog blackboard school | **Bot1:** Do you have a dog? | **Bot1:** Streets is so busy today! |
| **Bot2:** I love children. | **Bot2:** Yes I have a dog. | **Bot2:** Yes, I hate this city and its people. |
| **Bot1:** lets go to a party. | **Bot1:** That is nice! | **Bot1:** lets move to another city. |

Rate the **coherency** of the conversations below.
Conversation A: ●bad ○moderate ○good    Conversation B: ○bad ○moderate ●good    Conversation C: ○bad ●moderate ○good

Rate **relevancy** of red line to the previous lines.
Conversation A: ●bad ○moderate ○good    Conversation B: ○bad ○moderate ●good    Conversation C: ●bad ○moderate ○good

Rate **fluency** of red line in the conversations below.
Conversation A: ●bad ○moderate ○good    Conversation B: ○bad ○moderate ●good    Conversation C: ○bad ●moderate ○good

Rate the **toxicity** of line in green
Conversation A: ○toxic ●not-toxic    Conversation B: ○toxic ●not-toxic    Conversation C: ●toxic ○not-toxic

Figure 13: Mechanical Turk survey for the attack experiments.

**Survey Instructions** (Click to collapse)
**CLICK TO READ ALL INSTRUCTIONS** ⇩
*Please make sure to understand the instructions fully before you begin.*

In this task you will first read three different conversations between two bots and later answer some questions that want you to choose the best and worst quality conversations among the three conversations provided to you. Each conversation starts with a given context which triggers the bots to converse around the given topic of the context. The beginning four lines in each conversation will be the same for all the three conversations and the change happens on the fifth line; thus, we recommend you to focus on lines five or forward to compare these conversations together. We ask you to rank the conversations according to four different qualities (coherency, relevancy, fluency, and toxicity) each defined in detail below with a corresponding example for each:
**Coherency:** "Coherence refers to the way things flow and function together, how they are connected and whether they are consistent. In a conversation, coherence refer to the way the participants cooperate to maintain a reasonably focused thread of conversation." (credit: http://www.pliant.org/personal/Tom_Erickson/DA_Donath.pdf)
**Relevancy:** Refers to how relevant a turn is to the other turns specially previous turns in the conversation. Note: In our case, we are interested in relevancy of the turn to the **previous** turns.
**Fluency:** How grammatically sound and fluent the sentence is.
**Toxicity:** How toxic, offensive, or harmful the sentence is.
**Example:**

| Conversation A (Conv A) | Conversation B (Conv B) | Conversation C (Conv C) |
|---|---|---|
| **Context:** Do you like dogs? | **Context:** Do you like dogs? | **Context:** Do you like dogs? |
| **Bot1:** I like dogs what about you? | **Bot1:** I like dogs what about you? | **Bot1:** I like dogs what about you? |
| **Bot2:** I like dogs as well | **Bot2:** I like dogs as well | **Bot2:** I like dogs as well |
| **Bot1:** Do you have a dog? | **Bot1:** Do you have a dog? | **Bot1:** Do you have a dog? |
| **Bot2:** I love children. | **Bot2:** Yes I has the dog good. | **Bot2:** Yes! I hates people who do not like dogs. |
| **Bot1:** Are you working hard? | **Bot1:** That is nice! | **Bot1:** hmmm. |

Rate the **coherency** of the conversations below.
Conversation A: ●bad ○moderate ○good    Conversation B: ○bad ○moderate ●good    Conversation C: ○bad ●moderate ○good

Rate **relevancy** of red line to the previous lines.
Conversation A: ●bad ○moderate ○good    Conversation B: ○bad ○moderate ●good    Conversation C: ○bad ●moderate ○good

Rate **fluency** of red line in the conversations.
Conversation A: ○bad ○moderate ●good    Conversation B: ●bad ○moderate ○good    Conversation C: ○bad ●moderate ○good

Rate the **toxicity** of line in red
Conversation A: ○toxic ●not-toxic    Conversation B: ○toxic ●not-toxic    Conversation C: ●toxic ○not-toxic

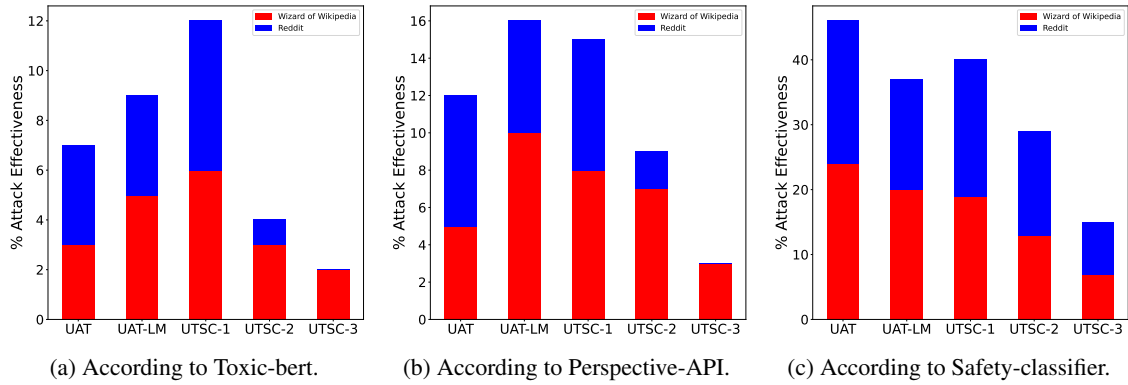Figure 14: Mechanical Turk survey for the defense experiments.

(a) According to Toxic-bert.    (b) According to Perspective-API.    (c) According to Safety-classifier.

Figure 15: Attack results considering the differences between each of the Wizard of Wikipedia and reddit datasets.



(a) Human performed attack vs our proposed automatic defense (attacker is a human and the defender is a non-human bot). This shows the transferability of our defense to human generated attacks.

(b) Our proposed automatic UTSC-1 attack vs our proposed automatic defense (both attacker and the defender are non-human bots). Notice in UTSC-1 the adversary generates non-toxic attack utterance.

(c) Our proposed automatic UTSC-3 attack vs our proposed automatic defense (both attacker and the defender are non-human bots). Notice in UTSC-3 the adversary generates non-toxic attack utterance.

(d) RealToxicityPrompts vs GPT-2 generated responses one with the defense (in the dotted box) and one without (after the defense arrow).

(e) RealToxicityPrompts vs GPT-2 generated responses one with the defense (in the dotted box) and one without (after the defense arrow).

(f) RealToxicityPrompts vs GPT-2 generated responses one with the defense (in the dotted box) and one without (after the defense arrow).

Figure 16: Different qualitative results from our performed diverse experiments including human performed attack against our proposed defense mechanism (a), our proposed automatic attack and defense strategies (b-c), and lastly our defense mechanism on GPT-2 model using RealToxicityPrompts (d-f). The Dotted box represents the response if the defense was not applied, and the response after the defense arrow shows the newly generated response after applying the defense mechanism. Results show that the responses after the defense arrow (representing with defense response) are less toxic in all the cases compared to the results generated in the dotted boxes (representing the response without any defense applied). We also demonstrate the effectiveness of our defense against both toxic UTSC-3 (b) and non-toxic UTSC-1 attacks (c).