

Outcome-Constrained Large Language Models for Countering Hate Speech

Anonymous EMNLP submission

Abstract

Counter speech (CS) that challenges or counteracts harmful or discriminatory messages is an effective way to diminish the influence of hate speech (HS). Automatic CS generation methods have been developed to assist efforts in combating online HS. Existing research focuses on generating CS with linguistic attributes, such as being polite, informative, and intent-driven. However, the real impact of CS in online environments is seldom considered. This study aims to develop methods for generating CS constrained by conversation outcomes and evaluate their effectiveness. We experiment with large language models (LLMs) to incorporate into the text generation process two desired conversation outcomes: low conversation incivility and non-hateful hater reentry. Specifically, we experiment with *instruction prompts*, *LLM finetuning*, and *LLM reinforcement learning (RL)*. Evaluation results show that our methods effectively steer the generation of conversational systems towards desired outcomes. Our analyses, however, show that there are differences in the quality and style of the generated CS.

1 Introduction

Hate speech (HS) has posed significant challenges to healthy and productive online communication. Counter speech (CS), which involves using constructive, positive, or factual responses to challenge or counteract HS, has shown to be effective in moderating online hostilities (Buerger, 2021), promoting productive user engagement (Miškolci et al., 2020), and educating online users (Blaya, 2019).

Automatic generation of CS has been researched to support moderators or individuals in their timely and effective efforts to fight HS. Synthetic CS datasets have been developed using crowdsourcing (Qian et al., 2019) and human-in-the-loop strategies (Chung et al., 2021). These datasets have promoted the development of CS generation mod-

els. However, the impact of CS in online environments has not been considered in the dataset creation. As a result, it is unknown whether generated CS elicits civil or hateful follow-up conversations.

Recent CS generation research focused on constrained generation with linguistic attributes (e.g., being polite, emotion-laden (Saha et al., 2022)), or embedded with knowledge (Chung et al., 2021). Questions about the impact of CS with such attributes linger. Previous research also found one of the barriers counterspeakers face is their inability to determine the potential impact of CS (Mun et al., 2024). However, there is a lack of research on generating outcome-oriented CS, e.g., speech that leads to desired outcomes such as de-escalating user conflicts or encouraging constructive engagement in follow-up conversations.

Notably, previous studies indicate that language may influence the development of a conversation, including discourse popularity (Horawalavithana et al., 2022), reentry behaviors (Wang et al., 2021), and the rise of hate speech (Liu et al., 2018). This leads to our research questions:

- How can constraints on conversation outcomes be incorporated into the development of LLMs for generating CS?
- How effective are these methods in generating outcome-oriented CS?

Unlike previous work that considers explicit linguistic attributes to guide language generation, we formulate CS generation to achieve desired outcomes (e.g., constructive user engagement). Our study holds potential for broader applications. Anticipating the direction of a conversation is crucial in crafting effective responses, allowing the conversation to meet the objectives of the interaction (e.g., reducing hate speech, altering user behavior, and promoting positive discourse). This study makes the following contributions: (i) introducing conversation outcomes as a constraint to guide the generation of CS, (ii) experimenting with LLMs

Prior Work	CS Constraint	HS Data	CS Generation Method
CONAN (Chung et al., 2019)	None	Islamophobic hate texts	Expert-based and LM data augmentation
Benchmark (Qian et al., 2019)	None	Reddit Gab	Crowdsourcing and LM generation
MultiCONAN (Fantón et al., 2021)	None	HS/CS from NGOs with multiple hate targets	LM generation with review/edits by experts
Knowledge (Chung et al., 2021)	Informative	CONAN	LM generation with information from knowledge repository
Generate-Prune (Zhu and Bhat, 2021)	Diverse and relevant	Benchmark CONAN	LM generation with quality classifier
COUNTERGEDİ (Saha et al., 2022)	Polite, detoxified, and emotional	Benchmark CONAN	DialoGPT and GEDİ for constraint generation
Intent (Gupta et al., 2023)	Multiple intents	CONAN MultiCONAN	QUARC with intent category representation and fusion
Ours	Expected outcomes	Benchmark CONAN, MultiCONAN	LLMs, LLM finetuning LLM RL

Table 1: Summary of recent work on counter speech generation, including dataset creation and modeling efforts.

for generating outcome-constrained CS using *instruction prompts*, *LLM finetuning*, and *LLM reinforcement learning (RL)*, and (iii) evaluating CS generation models with various metrics to understand the strengths and weaknesses of the methods.

2 Related Work

Generation of CS to HS Table 1 presents recent work on CS generation. Several CS datasets have been created. CONAN has CS written by NGO experts and augmented by language models (Chung et al., 2019); Benchmark was built with HS from Gab and Reddit and CS created by crowdsourcing workers (Qian et al., 2019); and MultiCONAN is a high-quality, high-quantity HS/CS dataset created by experts coupled with language model generation (Fantón et al., 2021). Several CS generation models have been built with these datasets (Halim et al., 2023; Tekiroğlu et al., 2020, 2022; Bonaldi et al., 2024) Unlike us, none of them consider the conversation outcomes of the generated CS.

Recently, researchers have investigated CS generation under selected constraints. Chung et al. (2021) proposed a CS generation pipeline grounded in external knowledge repositories to generate more informative and less biased replies. Zhu and Bhat (2021) generated more diverse and relevant CS and proposed a three-stage pipeline that generates CS candidates, prunes the ungrammatical ones, and selects the best instances. Saha et al. (2022) proposed an ensemble generative discriminator to generate more polite, detoxified, and emotion-laden CS. Gupta et al. (2023) developed IntentCONAN, where the generation of CS is conditioned on five

intents: informative, denouncing, question, positive, and humor. Similarly, Fraser et al. (2023) utilized ChatGPT to generate counter-stereotype text by incorporating countering strategies in queries. Hassan and Alikhani (2023) proposed prompting strategies based on discourse theories to generate more context-relevant CS. There are also studies on the generation of CS in languages other than English (e.g., Italian (Chung et al., 2020)). Unlike us, none of these previous works generate CS to elicit positive behaviors in the follow-up conversations.

Language Generation with Constraints Extensive studies have targeted language generation under complex lexical constraints such as formality (Jin et al., 2022), text with certain concepts (Lu et al., 2022), dialogue that takes latent variables (Bao et al., 2020), and knowledge-enhanced text (Yu et al., 2022a). Not all styles can be described explicitly as linguistic attributes. Indeed, some ‘styles’ can only be defined in a data-driven way based on the shared attributes across various datasets (Mou and Vechtomova, 2020). In this study, we generate CS very likely to lead to desired conversational outcomes.

Methods have been developed for constrained language generation. Wang and Wan (2018) proposed the SentiGAN framework to generate text with a given sentiment. Kumar et al. (2021) proposed MUCOCO to allow for controllable inference with multiple attributes as constraints to the optimization. Krause et al. (2021) developed GeDi, a discriminator-based approach to guide the decoding process in language generation. It enables

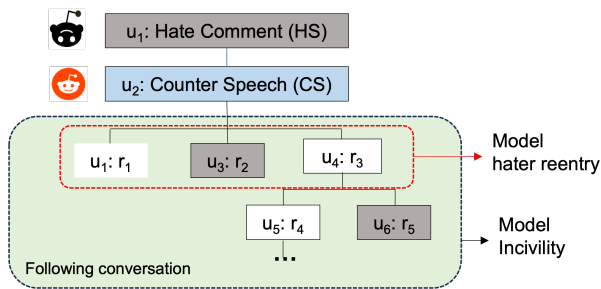


Figure 1: The two conversation outcomes to assess the conversation (green box) following up a counterspeech reply (blue box). Comments in the first layer of the conversation tree (i.e., direct replies) are used to model hater reentry. All comments in the conversation tree are used to model conversation incivility. Grey boxes indicate hateful comments; others are non-hateful.

text generation with desired or undesired attributes. Schick et al. (2021) proposed a self-debiasing approach to reduce the probability of language models generating problematic text. Unlike these previous efforts, we experiment with methods to adjust language model-generated texts to achieve specific conversational outcomes.

3 Methodology

3.1 Conversation Outcomes

Conversation outcomes refer to the result of a message in a conversation, which can be measured by the manner and characteristics of the follow-up conversations it elicits. According to previous studies, a combination of HS comment and its reply—regardless of whether it counters the hateful comment—can predict future conversation engagement and incivility (Liu et al., 2018; Yu et al., 2024). In this study, we explore two types of conversation outcome modeling: conversation incivility and hater reentry behavior (Figure 1). Based on the modeling results, we build conversation outcome classifiers that use the text of HS and CS comments to predict the incivility level or hater reentry type.

Conversation Incivility. Conversation incivility is a metric to measure the outcome based on the number of civil and uncivil comments as well as the unique authors involved in the discourse (Yu et al., 2024). Intuitively, the more uncivil (or less civil) the comments, the worse the outcome; uncivil comments from many authors are worse than those from just a few. Formally, it is defined as $S(r) = \alpha U(r) - (1 - \alpha)C(r)$, where $U(r)$ refers to uncivil behavior and $C(r)$ to civil behav-

ior. For each user i ($i = 0, 1, 2, \dots, k$), n_{ui} is defined as the number of uncivil comments by user i , and n_{ci} as the number of civil comments. Then, $U(r) = \sum_{i=1}^k \sqrt{n_{ui}}$ and $C(r) = \sum_{i=1}^k \sqrt{n_{ci}}$. α is used to adjust the weight of civil and uncivil behaviors. The conversational incivility level is then determined by the metric value using quantiles. Previous studies show that given two CS replies to a HS comment, models taking into account the text of the HS and CS comments accurately predict which of the two CS replies will lead to more civil follow-up conversations (Yu et al., 2024, binary classification, F1=0.66–0.75). We will use *civility* to refer to low conversation incivility, the desired outcome, in the remainder of the paper.

Hater Reentry Behavior. After a CS reply to a hate comment, the hate instigator may exhibit different behaviors. Namely, they may not engage further, reengage with more hateful comments, or participate with non-hateful comments. The outcome can be determined based on whether the following comments have one that is from the hater and whether this comment is hate speech. The non-hateful hater reentry is the most desirable, as it signals that the CS reply encouraged the individual to change his behavior (Baider, 2023). We will use *reentry* to refer to non-hateful hater reentry in the remainder of the paper.

3.2 Outcome-Constrained Counter Speech Generation

We explore the following methods to incorporate the outcome constraints into the generation process.

Instruction Prompts LLMs are capable of understanding natural conversations and generating replies. The straightforward strategy is to ask LLMs to generate replies considering the potential outcomes of the follow-up conversation. This explores whether LLMs might pick up information from the instruction and generate responses toward the desired outcomes. The prompts are as follows:

1. *Baseline*: No explicit expected outcomes.

```
User: "Here is a hate comment:
<Hate Comment>.
Please write a counter speech
reply to the hate comment."
```

2. *Civility*: Instruction with low conversation incivility as a desired outcome.

```
User: "Here is a hate comment:
<Hate Comment>."
```

231 Please write a counter speech
232 reply to the hate comment
233 so that it could lead to low
234 incivility in the follow-up
235 conversations."

236 **3. Reentry:** Instruction with non-hateful hater reen-
237 try as a desired outcome.

238 User: "Here is a hate comment:
239 <Hate Comment>.
240 Please write a counter speech
241 reply to the hate comment so
242 that the hater comes back and
243 has constructive engagement."

244 There are different ways to set these outcome-
245 constrained instructions, which might affect results.
246 We adopt the instructions above as baselines for
247 comparison purposes.

248 When given instructions, LLMs can generate one
249 or multiple CS replies. In addition to experimenting
250 with the first generated reply, we follow (Zhu and
251 Bhat, 2021) and also use a *Generate and Select*
252 method to generate multiple replies and select the
253 ones predicted to have desired outcomes according
254 to conversation outcomes classifiers (Section 3.1).

255 **LLM Finetuning** LLMs may not be fully
256 optimized for generating texts with specific
257 constraints—in our case, desired conversation out-
258 comes. The finetuning process can tailor LLMs to
259 learn the task of interest. To guide the LLM in gen-
260 erating outcome-constrained CS, we finetune the
261 model with datasets containing conversations with
262 the desired outcomes: the HS/CS pairs followed by
263 low conversation incivility (Yu et al., 2022b) and
264 the pairs that have non-hateful hater reentry. We
265 use the Parameter-Efficient Fine-Tuning (PEFT)
266 with Low-Rank Adaptation (LoRA) method (Hu
267 et al., 2021) to finetune LLMs.

268 **Reinforcement Learning with LLM (RL)** This
269 method integrates the conversation outcome clas-
270 sifiers (Section 3.1) as a reward function to guide
271 the training process, which includes three steps.
272 First, a hate comment is used as a query to get the
273 response generated by an LLM. The initial model
274 serves as a baseline for generating CS. Second, the
275 HS and generated response are fed into the conver-
276 sation outcomes classifiers to obtain their outcome
277 labels for assigning rewards. Specifically, pairs
278 with low incivility or non-hateful reentry will be re-
279 warding higher. Third, we maximize the probability
280 of the desired outcomes in the text generation pro-
281 cess. The HS/CS pairs are used to calculate the log

probabilities of tokens in the trained and the base
model. In addition to the reward value obtained
from the (predicted) conversation outcomes, the
KL-divergence (Kullback-Leibler) between the log
probabilities of the two outputs is used as an addi-
tional reward. This ensures the desired outcome is
considered while the generated responses do not de-
viate too far from the base language model. In sum-
mary, the reward is computed as $R = r - \beta * KL$.
We train the model with the Proximal Policy Opti-
mization (PPO) step until local stability is achieved.

3.3 Evaluation

Desired Conversation Outcome Metrics The eval-
uation aims to assess the ability of these methods to
generate CS that is more likely to achieve desired
outcomes. As it would be difficult—and arguably
unethical—to post the generated text to conversa-
tions on social media platforms to observe the real
outcomes, we adopt an approach that has been used
before (Saha et al., 2022; Tekiroğlu et al., 2022;
Halim et al., 2023; Gupta et al., 2023). That is,
we use the conversation incivility level classifier
and the hater reentry classifier (Section 3.1) trained
with real conversation data to make predictions
with the HS and generated CS pairs. Although the
accuracy of the classifiers is not perfect, given two
CS replies, these classifiers reliably identify the one
that will lead to better outcomes (Yu et al., 2024,
binary classification, F1=0.66–0.75). Thus, they
serve as a proxy to compare methods to generate
outcome-constrained CS. Additionally, we conduct
human assessments for reliability purposes.

Human Assessments Human assessment focuses
on three aspects: suitability, relevance, and effec-
tiveness. *Suitability* is measured by considering: (i)
whether the linguistic style of the reply suits the
conversation and (ii) whether the reply follows the
civil rules of the environment. *Relevance* evalu-
ates the appropriateness of the generated text with
respect to the content of the hate comment. *Effec-
tiveness* is evaluated based on whether it can stop
the spread of hate and foster constructive conversa-
tions, as perceived by human annotators. Two grad-
uate assistants, a male and female aged between 20
and 30, who are proficient in English and familiar
with social media, assist with the evaluation. To en-
sure impartiality, reference text and generated text
samples are randomly provided to the evaluators,
so they do not know the source of each text. The
agreement rate is calculated to assess reliability.

Stylistic Metrics The generated CS is evaluated by stylistic metrics commonly used in previous studies (Chung et al., 2021; Zhu and Bhat, 2021; Tekiroğlu et al., 2022). We calculate the similarity of generated CS against a reference dataset consisting of human-generated CS with the BLEU score (Chen and Cherry, 2014), ROUGE (Lin, 2004), METEOR (Banerjee and Lavie, 2005), and BERTScore (Zhang et al., 2019). The quality of generated texts is evaluated by the GRUEN metrics (Zhu and Bhat, 2020), including dimensions of grammaticality, redundancy, focus, and GRUEN score. The same scores are also calculated for the reference dataset for comparison purposes. Finally, we calculate the type-token ratio and distinct-n-grams to evaluate the diversity of generated texts (Fantón et al., 2021).

4 Experiments

4.1 Conversation Outcomes Classifiers

Data to Build Conversation Outcomes Classifiers we use Reddit data collected from 39 subreddits likely to contain abusive content (Vidgen et al., 2021). The hate comments are identified based on hate classifiers (Qian et al., 2019). Then, we collect replies to the hate comments and identify CS in the replies referring to Yu et al. (2022b). For each CS to the HS, we collect the follow-up replies and detect whether each one is hate speech based on hate classifiers. We use all follow-up replies to calculate the conversation incivility with $\alpha = 0.8$ and determine the incivility level by quantiles. The direct replies following CS are used to identify hater reentry behavior: whether the hate instigator reenters and the comment is non-hateful. Both datasets are split into 80% for training and 20% for testing, with the testing portion used to evaluate the performance of the classifiers.¹

Classification Model and Performance As this study is not aimed at the best performance in the classification tasks, we use the RoBERTa model (Liu et al., 2019) to train outcome classifiers. The texts of HS/CS pairs are used to predict the incivility level and the hater reentry behavior. The detailed classification results can be seen in Table 5 and 6 in A.4. Although the classification results are somewhat low, these suboptimal classifiers are enough to defeat the baseline and differentiate CS that will lead to high or low incivility in the follow-up conversation, as shown by (Yu et al., 2024). The

¹Data and models available at Github upon acceptance.

accuracy for identifying non-hateful reentry is the highest.

4.2 Generating Counter Speech

Dataset We use the benchmark-Reddit dataset (Qian et al., 2019) for CS generation and evaluation. The data contains HS from Reddit and CS generated by crowdsourcing workers. As we plan to explore the effect of this data in the finetuning and RL method, the HS/CS pairs are split randomly into 80% for training and 20% for testing and evaluation.

Instruction Prompts We use the Llama2-7b-chat model in our experiments to compare different methods, as we cannot train larger models like Llama2-13b-chat for *finetuning* and *RL* due to limited computing capacity. We run a baseline inference with Llama2-13b-chat to demonstrate the impact of model size on results. As the generation and evaluation are based on the benchmark-Reddit data, we apply the same system-level guideline: “Please generate a response in Reddit style” for all generations. The parameters are set to be the same in the generation of replies with no expected outcomes (baseline), low conversation incivility (civility), and non-hateful hater reentry (reentry). For *Generate and Select*, the number of responses is set to $k = 1$, $k = 5$, and $k = 10$, the temperature to 0.7, and the maximum length of reply to 512. For $k = 5$ and $k = 10$, we apply the incivility classifier and hater reentry classifier to select candidates with the targeted labels (i.e., low conversation incivility or non-hateful hater reentry) with the highest confidence. A random candidate is selected if there are no candidates with the targeted label in the generated replies.

Finetuning The Llama2-7b-chat model is finetuned with the HS and CS pairs that are followed with low conversation incivility or non-hateful reentry in the training data. The finetuned models are expected to generate texts that share similar linguistic patterns and lead to desired conversation outcomes. Additionally, we fine-tune models with reference datasets, including benchmark-Reddit, benchmark-Gab, CONAN, and MultiCONAN (see model details in A.2). This is to compare whether models developed with the existing CS datasets can generate CS with desired outcomes and the effects of these datasets on guiding CS generation.

Reinforcement Learning We use the Llama2-7b-chat as the base model for the RL process. The

reward for the RL process is generated based on the outcome classifiers: for the predicted categories of conversation incivility low, medium, and high, corresponding discrete rewards are assigned in descending order, namely 2, 1, and 0; for hater reentry classification, the reward for non-hateful reentry, no reentry, and hateful reentry is 2, 1, and 0, respectively. We also use the Llama-2-7b-chat finetuned with the benchmark-Reddit dataset. The finetuned model can learn the CS to HS patterns, therefore the model trained with RL can generate CS that has similar linguistic patterns with the CS in the benchmark-Reddit dataset while having a higher probability of leading to expected conversation outcomes. The hyperparameters are shown in A.2. We leave exploring RL with other finetuned models for future work.

5 Results and Analysis

All methods are evaluated with the same test set from the benchmark-Reddit. The Llama2-7b-chat sometimes avoids responding to queries the model determines to be inappropriate and generates empty responses. Table 2 shows the ratio of non-empty, noted as valid, responses by each model. Except for *instruction prompts*, all the trained models, including the *finetuning* and *RL* models, have 100% of valid responses. In *instruction prompts*, the valid response rate increases when using a more powerful model (Llama2-13b-chat), forcing the model to generate more candidates, or asking the model to generate CS with constrained queries.

Expected Outcomes In the task of generating texts with low conversation incivility, we observe the following insights: (i) The CS generated by a more powerful model (Llama2-13b-chat) has a higher proportion of samples leading to low incivility. (ii) Prompt queries with the constraint of low incivility can increase the probability of generating CS leading to the expected outcome. (iii) The *generate and select* strategy leads to more CS with the desired outcomes. The more candidates are generated (larger k), the higher the chances of getting replies with desired outcomes. (iv) The performance of *finetuning* methods in generating texts with expected outcomes is relatively inferior to others. (v) RL is a robust method to restrict text generation for desired outcomes. Both RL with Llama2-7b-chat and finetuned Llama2-7b-chat generate more responses with desired outcomes than the baseline models and *finetuning*. (vi) Human-

generated CS without consideration of outcomes in the benchmark-Reddit may fail to lead to expected conversation outcomes. only 760 samples (27%) are classified as having low conversation incivility.

The evaluation with the hater-reentry classifier further validates most insights. Larger models, prompts with desired outcomes, generate and select, and RL models generate more CS with desired outcomes.

Similarity to Reference Texts We evaluate the similarity of generated texts to the CS in the benchmark-Reddit data. We do not claim the CS to HS by the benchmark-Reddit data are standards. Instead, they serve as a baseline for us to understand whether the LLM-generated texts are different from human-generated ones and how different. We calculate multiple similarity metrics. Results show the metrics are highly correlated (Table 9 in the A.5). Hence, we only present the results of METEOR and BERTScore in Table 2.

All the METEOR values are low, with the average values ranging from 0.06 to 0.14. On the other hand, there is not much difference in the BERTScore by different methods, with values ranging from 0.80 to 0.86. The difference between METEOR and BERTScores indicates that LLM-generated replies have high semantic similarity to reference CS, but the wording used in LLM-generated texts is different. Notably, even without finetuning or RL, LLMs are still capable of generating CS with similar meanings to reference texts (baseline generation BERTScore 0.8).

Quality of Generated Texts Table 3 presents the evaluation using stylistic metrics. Grammaticality scores measure grammatical correctness. Texts generated by language models generally have higher grammatical scores than the reference (0.77), except the ones finetuned with Reddit conversation data: civility (0.77) and reentry (0.76). These finetuned models might have learned informal expressions on social media, thus they generate CS with a lower grammaticality score. Texts generated by LLMs without finetuning or RL have more redundancy indicated by lower redundancy scores. After adding expected outcomes as constraints, LLM-generated CS contains less redundancy. The focus scores of texts generated by *instruction prompts* are also much lower. In models with *finetuning* and *RL*, the focus scores are much higher.

Overall, texts generated by *finetuning* and *RL* have higher quality, reflected in dimensions of

Instruction Prompts	Valid (%)	Desired Outcomes		Similarity	
		Civility (%)	Reentry (%)	METEOR	BERTScore
Generate one based on (k=1)					
Baseline	83%	23%	18%	0.07 (0.08)	0.80 (0.03)
Baseline(13B)	94%	27%	35%	0.12 (0.07)	0.81 (0.04)
Civility	92%	54%	49%	0.12 (0.05)	0.83 (0.02)
Reentry	94%	44%	45%	0.12 (0.06)	0.82 (0.02)
Generate and select (k=5)					
p=baseline, c=civility	84%	55%	32%	0.10 (0.07)	0.81 (0.03)
p=baseline, c=reentry	85%	34%	49%	0.11 (0.07)	0.82 (0.03)
p=civility, c=civility	92%	81%	53%	0.12 (0.05)	0.82 (0.02)
p=reentry, c=reentry	92%	49%	83%	0.13 (0.05)	0.83 (0.01)
Generate and select (k=10)					
p=baseline, c=civility	87%	69%	36%	0.11 (0.07)	0.82 (0.02)
p=baseline, c=reentry	86%	41%	61%	0.11 (0.07)	0.82 (0.02)
p=civility, c=civility	92%	86%	55%	0.12 (0.05)	0.82 (0.02)
p=reentry, c=reentry	92%	50%	86%	0.13 (0.05)	0.83 (0.01)
Finetuning with CS Corpora					
CONAN	100%	23%	48%	0.09 (0.06)	0.85 (0.02)
MultiCONAN	100%	22%	48%	0.11 (0.06)	0.85 (0.02)
Benchmark-Gab	100%	10%	43%	0.12 (0.10)	0.86 (0.02)
Benchmark-Reddit	100%	11%	42%	0.13 (0.11)	0.86 (0.02)
Ours, with conversation outcomes					
Reddit-CS-civility	100%	18%	35%	0.08 (0.05)	0.84 (0.02)
Reddit-CS-reentry	100%	19%	35%	0.08 (0.05)	0.84 (0.02)
Reinforcement Learning (RL)					
Civility	100%	77%	71%	0.14 (0.05)	0.83 (0.01)
Reentry	100%	67%	62%	0.14 (0.05)	0.83 (0.01)
RL with Benchmark-Reddit finetuned LLM					
Civility	100%	30%	48%	0.13 (0.13)	0.85 (0.02)
Reentry	100%	18%	57%	0.07 (0.06)	0.86 (0.01)
Reference					
Benchmark-Reddit	100%	27%	37%	1.00 (0.00)	1.00 (0.00)

Table 2: Evaluation of (a) Desired Outcomes and (b) Similarity to the reference CS in Benchmark-Reddit. METEOR and BERTScore are calculated per sample. Mean (SD) is reported. *Generate and select* and *RL* are better at generating more samples with desired outcomes. Although the wording differs from the Reference CS (METEOR), the semantic relevance (BERTScore) is consistently high.

grammaticality, redundancy, focus, and the overall GRUEN score. In particular, the highest GRUEN scores are achieved by *RL* models.

Diversity and Novelty The three diversity metrics (i.e., TTR, number of unique unigrams, and number of unique bigrams) are highly correlated (Table 8 in A.5). TTR and the novelty metric (i.e., number of new unigrams) are presented in Table 3.

The TTR of generated texts significantly decreases when using expected outcomes in *instruction prompts* and *RL*. LLMs finetuned with appropriate datasets generate mode-diverse CS. The highest TTRs are achieved by LLM models finetuned with real Reddit conversation data, which usually contains diverse, informal expressions.

The novelty of generated texts is higher when conversation outcomes are considered in the generation. The number of new unigrams generated by untrained LLMs in the *instruction prompt* method is substantially higher than trained models with *finetuning* and *RL*.

Human Evaluation We choose generated texts constrained with low conversation incivility for human evaluation. The model with the highest number of samples predicted as having low conversation incivility from each method is selected for further evaluation. Hence, we randomly select 50 pairs of HS and generated CS from the *instruction prompts* with $p = civility$, $k = 10$, and $c = civility$, *finetuning* with CONAN, and *RL* with low incivility, respectively. Then, we mix the samples and ask annotators to label yes or no to the suitability, relevance, and effectiveness. The percentages of agreement for initial evaluation are 0.78, 0.92, and 0.64 separately for suitability, quality, and effectiveness. For the samples without an agreement, the annotators discuss and finalize an agreed annotation. Table 4 presents the evaluation results. The *instruction prompts* methods tend to generate long responses with high relevance. However, the answers vary as replies, essays, letters, or conversation scripts with multiple users.

	Text Quality			Diversity	Novelty	
	Grammaticality	Focus	Redundancy	GRUEN	TTR	New Tokens
Instruction Prompts						
Generate one based on						
Baseline	0.73 (0.10)	-0.05 (0.05)	-1.14 (12.56)	0.60 (0.18)	0.06	5384
Baseline (13B)	0.80 (0.07)	-0.09 (0.03)	-1.33 (23.22)	0.60 (0.21)	0.06	9231
Civility	0.84 (0.04)	-0.10 (0.01)	-0.19 (0.56)	0.61 (0.22)	0.03	7019
Reentry	0.83 (0.07)	-0.10 (0.02)	-0.11 (0.39)	0.64 (0.18)	0.03	6407
Generate and select k=5						
p=baseline, c=civility	0.78 (0.10)	-0.08 (0.04)	-0.33 (4.37)	0.62 (0.19)	0.06	7220
p=baseline, c=reentry	0.78 (0.10)	-0.08 (0.04)	-0.34 (6.42)	0.63 (0.18)	0.05	6794
p=civility, c=civility	0.84 (0.03)	-0.10 (0.01)	-0.23 (2.35)	0.59 (0.23)	0.04	7668
p=reentry, c=reentry	0.84 (0.02)	-0.10 (0.00)	-0.07 (0.21)	0.68 (0.12)	0.03	5224
Generate and select k=10						
p=baseline, c=civility	0.79 (0.09)	-0.08 (0.04)	-0.27 (2.27)	0.62 (0.20)	0.06	8000
p=baseline, c=reentry	0.80 (0.09)	-0.08 (0.04)	-0.20 (2.02)	0.64 (0.18)	0.05	6908
p=civility, c=civility	0.84 (0.03)	-0.10 (0.00)	-0.23 (0.48)	0.57 (0.24)	0.04	8024
p=reentry, c=reentry	0.84 (0.02)	-0.10 (0.00)	-0.06 (0.12)	0.68 (0.11)	0.03	5198
Finetuning with CS Corpora						
CONAN	0.81 (0.09)	-0.02 (0.04)	0.00 (0.03)	0.78 (0.11)	0.11	1982
MultiCONAN	0.83 (0.07)	-0.05 (0.05)	-0.12 (2.93)	0.76 (0.13)	0.09	2448
Benchmark-Gab	0.85 (0.06)	-0.01 (0.03)	0.00 (0.00)	0.83 (0.08)	0.02	111
Benchmark-Reddit	0.80 (0.09)	-0.04 (0.05)	0.00 (0.01)	0.77 (0.12)	0.03	147
Ours, with conversation outcomes						
Reddit-CS-civility	0.78 (0.09)	-0.04 (0.05)	-0.70 (7.78)	0.71 (0.17)	0.12	2858
Reddit-CS-reentry	0.78 (0.09)	-0.04 (0.05)	-0.70 (7.56)	0.71 (0.17)	0.11	2643
Reinforcement Learning (RL)						
Civility	0.85 (0.03)	-0.10 (0.00)	-0.04 (0.12)	0.71 (0.11)	0.03	5575
Reentry	0.84 (0.04)	-0.10 (0.00)	-0.06 (0.18)	0.69 (0.13)	0.03	6574
RL with Benchmark-Reddit finetuned LLM						
Civility	0.80 (0.02)	0.00 (0.00)	0.00 (0.00)	0.80 (0.02)	0.00	0
Reentry	0.87 (0.03)	0.00 (0.00)	0.00 (0.00)	0.87 (0.03)	0.01	12
Reference						
Benchmark-Reddit	0.77 (0.12)	-0.03 (0.05)	0.00 (0.01)	0.74 (0.13)	0.09	0

Table 3: Evaluation of Stylistic Metrics: Quality and Diversity. GRUEN and BERTScore are calculated per sample. Mean (SD) are reported. The quality of CS by *Instruction prompts* is relatively low. *LLM finetuning* with Reddit-CS generate texts with high diversity. *RL* with finetuned LLMs generate texts with reduced novelty.

Method	Suitability	Relevance	Effectiveness
Prompt	0.50	0.88	0.54
Finetuning	0.80	0.68	0.80
RL	0.74	0.76	0.72

Table 4: Proportion of samples labeled as yes for each evaluation dimension by methods.

Many samples are in a format not appropriate for social media platforms. Although the desired outcome metric shows *finetuning* is relatively inferior to other methods, the human evaluation shows the generated CS by *finetuning* and *RL* are usually suitable, and effective. It deserves further investigation into the reasons that explain the differences in desired outcome metrics and human assessment.

6 Conclusions

We present an initial exploration of methods for constrained generation of CS controlled by potential conversation outcomes. We incorporate the desired outcomes (i.e., low conversation incivility

and non-hateful hater reentry) into the text generation process through three methods: *instruction prompts*, *LLM finetuning*, and *LLM RL*. The text generation results are evaluated with desired conversation metrics, stylistic metrics, and human assessment. Results show that *instruction prompts* and *RL* generate CS with a higher probability of eliciting desired outcomes based on the prediction of outcome classifiers, while *finetuning* and *RL* generate more effective CS based on human assessments. The LLMs-generated texts consistently show high relevance to HS, but the wording differs.

The generated texts present different characteristics. CS generated by LLM without further training tends to be long, not suitable for the conversation context on social media, and with low quality based on GRUEN metrics and human assessment. Both *finetuning* and *RL* models generate CS with high quality with styles suitable for social media platforms. The experiments present different methods' strengths and weaknesses, enabling stakeholders to choose methods appropriate for their needs.

609 **Limitations**

610 The conversation outcome classifiers are not perfect as the texts of hate comments and replies only
611 partially contribute to the conversation outcomes.
612 Other influencing factors include the context of
613 the conversation and users' positions and identities.
614 While the outcome classifiers provide a convenient method for evaluation, they may introduce
615 bias into the evaluation process. Therefore, interpretations and conclusions drawn from these evaluations
616 should be considered with caution. Future work will explore more accurate and unbiased
617 classifiers to enhance text generation and evaluation. We use computing-based metrics for evaluating
618 similarity, quality of text, diversity, and novelty. Although these metrics are widely used, they
619 may present bias. More sophisticated evaluation methods and comprehensive human assessments
620 are needed to fully capture the multidimensional quality of the generated text. Text generation is
621 influenced by numerous factors, including the formulation of prompt queries, settings of LLMs for
622 text generation, fine-tuning language models with different datasets, variations in fine-tuning and
623 reinforcement learning settings, and size of language models. Further experiments are needed to
624 better understand the impact of these factors on text generation. The outcome classifiers are based on
625 Reddit conversation data, which may not transfer to other platforms. Experiments with different data
626 are to be done to understand communication patterns across platforms and the guiding effect of
627 cross-domain data.

642 **Ethics Statement**

643 The study has been through careful consideration of benefits and risks. First, we used data from
644 Reddit, which is considered a public space. Users consent to make their data available to third parties.
645 Second, user names and identities are encrypted to avoid the identification of users. Third, student
646 collaborators working on the data have been warned of the potential hateful content and are encouraged
647 to stop their work at any time. Fourth, the data will be shared for research purposes only. Although
648 releasing the dataset may raise risks, we believe the benefits of contributing to effective methods to
649 counter online hate outweighs the potential risks. Finally, the models developed may not be directly
650 applicable to the generation of CS to online hate. Instead, they could serve as valuable tools to assist

content moderation in crafting CS. Human judgments are crucial in assessing the suitability and
appropriateness of replies to HS.

662 **References**

- 663 Fabienne Baider. 2023. Accountability issues, online covert hate speech, and the efficacy of counter-speech. *Politics and Governance*, 11(2):249–260. 664 665
- 666 Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72. 667 668 669 670 671
- 672 Siqi Bao, Huang He, Fan Wang, Hua Wu, and Haifeng Wang. 2020. Plato: Pre-trained dialogue generation model with discrete latent variable. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 85–96. 673 674 675 676
- 677 Catherine Blaya. 2019. Cyberhate: A review and content analysis of intervention strategies. *Aggression and violent behavior*, 45:163–172. 678 679
- 680 Helena Bonaldi, Yi-Ling Chung, Gavin Abercrombie, and Marco Guerini. 2024. Nlp for counterspeech against hate: A survey and how-to guide. *arXiv preprint arXiv:2403.20103*. 681 682 683
- 684 Catherine Buerger. 2021. # iamhere: Collective counterspeech and the quest to improve online discourse. *Social Media+ Society*, 7(4):20563051211063843. 685 686
- 687 Boxing Chen and Colin Cherry. 2014. A systematic comparison of smoothing techniques for sentence-level bleu. In *Proceedings of the ninth workshop on statistical machine translation*, pages 362–367. 688 689 690
- 691 Yi-Ling Chung, Elizaveta Kuzmenko, Serra Sinem Tekiroğlu, and Marco Guerini. 2019. Conan-counter narratives through nichesourcing: a multilingual dataset of responses to fight online hate speech. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2819–2829. 692 693 694 695 696 697
- 698 Yi-Ling Chung, Serra Sinem Tekiroğlu, and Marco Guerini. 2020. Italian counter narrative generation to fight online hate speech. In *Proceedings of the Seventh Italian Conference on Computational Linguistics (CLIC-it 2020)*, volume 2769. 699 700 701 702
- 703 Yi-Ling Chung, Serra Sinem Tekiroğlu, and Marco Guerini. 2021. Towards knowledge-grounded counter narrative generation for hate speech. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 899–914. 704 705 706 707
- 708 Margherita Fanton, Helena Bonaldi, Serra Sinem Tekiroğlu, and Marco Guerini. 2021. Human-in-the-loop for data collection: a multi-target counter narrative dataset to fight online hate speech. In *Proceedings of the 59th Annual Meeting of the Association for* 709 710 711 712

713			
714		<i>Computational Linguistics and the 11th International</i>	
715		<i>Joint Conference on Natural Language Processing</i>	
		(Volume 1: Long Papers), pages 3226–3240.	
716	Kathleen C Fraser, Svetlana Kiritchenko, Isar Ne-		
717	jadgholi, and Anna Kerkhof. 2023. What makes		
718	a good counter-stereotype? evaluating strategies for		
719	automated responses to stereotypical text. In <i>Pro-</i>		
720	<i>ceedings of the First Workshop on Social Influence</i>		
721	<i>in Conversations (SICon 2023)</i> , pages 25–38.		
722	Rishabh Gupta, Shaily Desai, Manvi Goel, Anil		
723	Bandhakavi, Tanmoy Chakraborty, and Md Shad		
724	Akhtar. 2023. Counterspeeches up my sleeve! in-		
725	tent distribution learning and persistent fusion for		
726	intent-conditioned counterspeech generation. <i>arXiv</i>		
727	<i>preprint arXiv:2305.13776</i> .		
728	Sadaf MD Halim, Saquib Irtiza, Yibo Hu, Latifur Khan,		
729	and Bhavani Thuraisingham. 2023. Wokegpt: Im-		
730	proving counterspeech generation against online hate		
731	speech by intelligently augmenting datasets using a		
732	novel metric. In <i>2023 International Joint Conference</i>		
733	<i>on Neural Networks (IJCNN)</i> , pages 1–10. IEEE.		
734	Sabit Hassan and Malihe Alikhani. 2023. Discgen:		
735	A framework for discourse-informed counterspeech		
736	generation. In <i>Proceedings of the 13th International</i>		
737	<i>Joint Conference on Natural Language Processing</i>		
738	<i>and the 3rd Conference of the Asia-Pacific Chapter</i>		
739	<i>of the Association for Computational Linguistics</i>		
740	(Volume 1: Long Papers), pages 420–429.		
741	Sameera Horawalavithana, Nazim Choudhury, John		
742	Skvoretz, and Adriana Iamnitchi. 2022. Online dis-		
743	cussion threads as conversation pools: predicting the		
744	growth of discussion threads on reddit. <i>Computa-</i>		
745	<i>tional and Mathematical Organization Theory</i> , pages		
746	1–29.		
747	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu,		
748	Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,		
749	et al. 2021. Lora: Low-rank adaptation of large lan-		
750	guage models. In <i>International Conference on Learn-</i>		
751	<i>ing Representations</i> .		
752	Di Jin, Zhijing Jin, Zhiting Hu, Olga Vechtomova,		
753	and Rada Mihalcea. 2022. Deep learning for text		
754	style transfer: A survey. <i>Computational Linguistics</i> ,		
755	48(1):155–205.		
756	Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann,		
757	Nitish Shirish Keskar, Shafiq Joty, Richard Socher,		
758	and Nazneen Fatema Rajani. 2021. Gedi: Genera-		
759	tive discriminator guided sequence generation. In		
760	<i>Findings of the Association for Computational Lin-</i>		
761	<i>guistics: EMNLP 2021</i> , pages 4929–4952.		
762	Sachin Kumar, Eric Malmi, Aliaksei Severyn, and Yu-		
763	lia Tsvetkov. 2021. Controlled text generation as		
764	continuous optimization with multiple constraints.		
765	<i>Advances in Neural Information Processing Systems</i> ,		
766	34:14542–14554.		
767	Chin-Yew Lin. 2004. Rouge: A package for automatic		
768	evaluation of summaries. In <i>Text summarization</i>		
769	<i>branches out</i> , pages 74–81.		
	Ping Liu, Joshua Guberman, Libby Hemphill, and Aron		770
	Culotta. 2018. Forecasting the presence and intensity		771
	of hostility on instagram using linguistic and social		772
	features. In <i>Proceedings of the International AAAI</i>		773
	<i>Conference on Web and Social Media</i> , volume 12.		774
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-		775
	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,		776
	Luke Zettlemoyer, and Veselin Stoyanov. 2019.		777
	Roberta: A robustly optimized bert pretraining ap-		778
	proach. <i>arXiv preprint arXiv:1907.11692</i> .		779
	Ximing Lu, Sean Welleck, Peter West, Liwei Jiang,		780
	Jungo Kasai, Daniel Khashabi, Ronan Le Bras, Lian-		781
	hui Qin, Youngjae Yu, Rowan Zellers, et al. 2022.		782
	Neurologic a* esque decoding: Constrained text gen-		783
	eration with lookahead heuristics. In <i>Proceedings of</i>		784
	<i>the 2022 Conference of the North American Chap-</i>		785
	<i>ter of the Association for Computational Linguistics:</i>		786
	<i>Human Language Technologies</i> , pages 780–799.		787
	Jozef Miškolci, Lucia Kováčová, and Edita Rigová.		788
	2020. Countering hate speech on facebook: The		789
	case of the roma minority in slovakia. <i>Social Science</i>		790
	<i>Computer Review</i> , 38(2):128–146.		791
	Lili Mou and Olga Vechtomova. 2020. Stylized text gen-		792
	eration: Approaches and applications. In <i>Proceed-</i>		793
	<i>ings of the 58th Annual Meeting of the Association</i>		794
	<i>for Computational Linguistics: Tutorial Abstracts</i> ,		795
	pages 19–22.		796
	Jimin Mun, Cathy Buerger, Jenny T Liang, Joshua Gar-		797
	land, and Maarten Sap. 2024. Counterspeakers’ per-		798
	spectives: Unveiling barriers and ai needs in the fight		799
	against online hate. In <i>Proceedings of the CHI Con-</i>		800
	<i>ference on Human Factors in Computing Systems</i> ,		801
	pages 1–22.		802
	Jing Qian, Anna Bethke, Yinyin Liu, Elizabeth Belding,		803
	and William Yang Wang. 2019. A benchmark dataset		804
	for learning to intervene in online hate speech. In		805
	<i>Proceedings of the 2019 Conference on Empirical</i>		806
	<i>Methods in Natural Language Processing and the 9th</i>		807
	<i>International Joint Conference on Natural Language</i>		808
	<i>Processing (EMNLP-IJCNLP)</i> , pages 4755–4764.		809
	Punyajoy Saha, Kanishk Singh, Adarsh Kumar, Binny		810
	Mathew, and Animesh Mukherjee. 2022. Coun-		811
	tergedi: A controllable approach to generate po-		812
	lite, detoxified and emotional counterspeech. <i>arXiv</i>		813
	<i>preprint arXiv:2205.04304</i> .		814
	Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021.		815
	Self-diagnosis and self-debiasing: A proposal for re-		816
	ducing corpus-based bias in nlp. <i>Transactions of the</i>		817
	<i>Association for Computational Linguistics</i> , 9:1408–		818
	1424.		819
	Serra Sinem Tekiroğlu, Helena Bonaldi, Margherita		820
	Fanton, and Marco Guerini. 2022. Using pre-trained		821
	language models for producing counter narratives		822
	against hate speech: a comparative study. In <i>Find-</i>		823
	<i>ings of the Association for Computational Linguistics:</i>		824
	<i>ACL 2022</i> , pages 3099–3114.		825

826 Serra Sinem Tekiroğlu, Yi-Ling Chung, and Marco
827 Guerini. 2020. Generating counter narratives against
828 online hate speech: Data and strategies. In *Proceed-*
829 *ings of the 58th Annual Meeting of the Association*
830 *for Computational Linguistics*, pages 1177–1190.

831 Bertie Vidgen, Dong Nguyen, Helen Margetts, Patricia
832 Rossini, and Rebekah Tromble. 2021. Introducing
833 cad: the contextual abuse dataset. In *Proceedings of*
834 *the 2021 Conference of the North American Chap-*
835 *ter of the Association for Computational Linguistics:*
836 *Human Language Technologies*, pages 2289–2303.

837 Ke Wang and Xiaojun Wan. 2018. Sentigan: Generating
838 sentimental texts via mixture adversarial networks.
839 In *Proceedings of the Twenty-Seventh International*
840 *Joint Conference on Artificial Intelligence (IJCAI-*
841 *18)*, pages 4446–4452.

842 Lingzhi Wang, Xingshan Zeng, Huang Hu, Kam-Fai
843 Wong, and Daxin Jiang. 2021. Re-entry prediction
844 for online conversations via self-supervised learning.
845 In *Findings of the Association for Computational*
846 *Linguistics: EMNLP 2021*, pages 2127–2137.

847 Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu,
848 Qingyun Wang, Heng Ji, and Meng Jiang. 2022a. A
849 survey of knowledge-enhanced text generation. *ACM*
850 *Computing Surveys*, 54(11s):1–38.

851 Xinchun Yu, Eduardo Blanco, and Lingzi Hong. 2022b.
852 Hate speech and counter speech detection: Conversa-
853 tional context does matter. In *Proceedings of the*
854 *2022 Conference of the North American Chapter of*
855 *the Association for Computational Linguistics: Hu-*
856 *man Language Technologies*, pages 5918–5930.

857 Xinchun Yu, Eduardo Blanco, and Lingzi Hong. 2024.
858 Hate cannot drive out hate: Forecasting conversation
859 incivility following replies to hate speech. In *Pro-*
860 *ceedings of the International AAAI Conference on*
861 *Web and Social Media*, volume 18, pages 1740–1752.

862 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Wein-
863 berger, and Yoav Artzi. 2019. Bertscore: Evaluating
864 text generation with bert. In *International Confer-*
865 *ence on Learning Representations*.

866 Wanzheng Zhu and Suma Bhat. 2020. Gruen for evalu-
867 ating linguistic quality of generated text. In *Find-*
868 *ings of the Association for Computational Linguistics:*
869 *EMNLP 2020*, pages 94–108.

870 Wanzheng Zhu and Suma Bhat. 2021. Generate, prune,
871 select: A pipeline for counterspeech generation
872 against online hate speech. In *Findings of the Associ-*
873 *ation for Computational Linguistics: ACL-IJCNLP*
874 *2021*, pages 134–149.

875 A Appendices

876 A.1 Computing Resources

877 The computational resources used in this research
878 include a high-performance server equipped with
879 three Quadro RTX 8000 GPUs, 128G memory, and
880 a 4T disk.

881 A.2 Hyperparameters

882 LLM Finetuning: We use PEFT LoRA for the
883 finetuning process. The LoRA configuration has
884 $r = 16$, $alpha = 32$, $dropout = 0.05$, and bias is
885 “none”. The hyperparameters are as follows: the
886 learning rate is $1e-4$, the number of epochs is 1,
887 and the warmup ratio is 0.1.

888 LLM RL: The reward trainer uses the RoBERTa
889 base model, the learning rate is $1e-5$, the batch size
890 is 16, and the number of epochs is 5. In the PPO
891 process, the generation component has $top_k = 0$,
892 $top_p = 1.0$, $do_sample = True$, and the max
893 length is 256. The PPO configuration has a learning
894 rate of $1.41e-5$, a batch size of 32, and an initial
895 KL coefficient of 0.1.

896 A.3 Dataset License and Use

897 The Benchmark dataset by Qian et al. (2019)
898 is under the Creative Commons Attribution-
899 NonCommercial 4.0 International Public License.
900 The CONAN and MultiCONAN datasets can
901 be used for research purposes with proper cita-
902 tion (Chung et al., 2019; Fanton et al., 2021). The
903 benchmark-Reddit data contains 5,020 unique con-
904 versations with hate speech identified. Each hate
905 speech comment has multiple responses. We ex-
906 tracted the hate speech from conversations and their
907 CS responses, generating 14,208 valid HS/CS pairs,
908 noted as the benchmark-Reddit data. The testing
909 data includes 2,843 pairs of HS/CS.

910 A.4 Evaluation Results of Conversation 911 Outcome Classifiers

912 Table 5 presents the evaluation of the conversa-
913 tion incivility classifier. The baseline is calculated
914 assuming all test samples are assigned with the
915 majority label, Medium. Although the classifica-
916 tion results are somewhat low, these suboptimal
917 classifiers are enough to defeat the baseline and
918 differentiate CS that will lead to high or low in-
919 civility in the follow-up conversation (Yu et al.,
920 2024, binary classification, $F1=0.66-0.75$). Table 6
921 presents the evaluation of the hater reentry classi-
922 fier. The baseline is calculated assuming all test
923 samples are assigned with the majority label, non-
924 hateful reentry. The non-hateful reentry class has
925 the highest F1 of 0.61.

926 A.5 Evaluation Metrics

927 Table 7 shows the number of samples in each class
928 based on the prediction of the conversation incivil-

	High			Medium			Low			Weighted Average		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Baseline	0.00	0.00	0.00	0.49	1.00	0.66	0.00	0.00	0.00	0.24	0.49	0.32
Incivility	0.43	0.32	0.36	0.55	0.66	0.60	0.32	0.27	0.29	0.46	0.48	0.46

Table 5: Evaluation results of the conversation incivility classifier.

	Hate reentry			No reentry			Non-hate reentry			Weighted Average		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Baseline	0.00	0.00	0.00	0.00	0.00	0.00	0.49	1.00	0.66	0.16	0.33	0.22
Reentry	0.32	0.20	0.25	0.52	0.41	0.46	0.54	0.70	0.61	0.49	0.51	0.46

Table 6: Evaluation results of the hater reentry classifier.

ity classifier and the hate re-entry classifier.

Table 8 presents the correlation coefficients between diversity metrics (i.e., type-token ratio, distinct-1, and distinct-2) and novelty metrics (i.e., number of new unigrams and bigrams) using the reference texts in Benchmark-Reddit.

Table 9 presents the correlation of metrics that evaluate the relevance of generated texts to reference texts in Benchmark-Reddit.

Table 10 presents a relatively good and bad examples of generated texts by different methods². Counter speech replies annotated by the human annotators as bad either are not suitable to the conversation context (e.g., example(2)), not a counter speech (e.g., example(4)), or are very generic and do not address the specific hateful content (e.g., example(6)).

A.6 AI Use

We acknowledge the use of code-writing assistance GitHub Copilot. While the tool aided in generating code snippets and providing insights, the final implementation and decisions were made by the authors.

²The examples in this paper contain hateful content. We cannot avoid it due to the nature of our work.

Category	Model	Conversation Incivility			Hater Reentry		
		High	Medium	Low	No reentry	Hateful	Non-hateful
Generation	baseline	291	1733	652	1422	748	506
	baseline(13B)	686	1214	776	752	937	987
	civility	412	657	1547	876	346	1394
	reentry	629	794	1253	910	476	1290
Prompt and Select	p=baseline k=5 c=civility	195	855	1566	1117	595	904
	p=civility k=5 c=civility	134	176	2306	849	253	1514
	p=baseline k=5 c=reentry	415	1240	961	771	443	1402
	p=reentry k=5 c=reentry	914	312	1390	64	186	2366
	p=baseline k=10 c=civility	114	537	1965	1070	511	1035
	p=civility k=10 c=civility	73	100	2443	828	222	1566
	p=baseline k=10 c=reentry	444	994	1178	511	371	1734
	p=reentry k=10 c=reentry	890	295	1431	25	160	2431
LLM Finetune	civility	953	1298	592	881	954	1008
	reentry	939	1417	487	731	1152	960
	CONAN	1429	752	662	438	1031	1374
	MultiCONAN	1386	835	622	559	931	1353
	Benchmark-Reddit	1775	757	311	510	1149	1184
	Benchmark-Gab	1974	585	284	533	1076	1234
LLM TRL	civility	239	423	2181	292	540	2011
	reentry	481	461	1901	408	661	1774
	bm_reddit_ft_civility	66	1917	860	448	1036	1359
	bm_reddit_ft_reentry	1212	1130	501	222	992	1629
Reference	benchmark_reddit	1245	838	760	683	1117	1043

Table 7: Evaluation results of conversation incivility and hater reentry classifiers.

	TTR	distinct-1	distinct-2	#new_unigram	#new_bigram
TTR	1	0.990	0.971	-0.219	-0.298
distinct-1		1	0.972	-0.287	-0.364
distinct-2			1	-0.086	-0.160
#new_unigram				1	0.980
#new_bigram					1

Table 8: Correlation coefficients of diversity metrics. TTR (type-token ratio) highly correlates with distinct-1 and distinct-2. Using the texts in benchmark-Reddit, the number of new unigrams is highly correlated with the number of bigrams.

	BLEU	ROUGE-r	ROUGE-p	ROUGE-f	METEOR	BERT-p	BERT-r	BERT-f
BLEU	1	0.799	0.997	0.990	0.973	0.870	0.983	0.952
ROUGE-r		1	0.829	0.874	0.916	0.466	0.866	0.649
ROUGE-p			1	0.996	0.983	0.848	0.985	0.939
ROUGE-f				1	0.995	0.806	0.99	0.914
METEOR					1	0.759	0.988	0.884
BERT-p						1	0.835	0.974
BERT-r							1	0.938
BERT-f								1

Table 9: Correlation coefficients of similarity metrics. Almost all metrics are highly correlated.

ID	Method	Evaluation	Hate Comment	Generated Reply
1	Instruction Prompts	Good	Wow. Can't believe anyone would actually vote for that cunt.	I understand that many people are upset about the election results, but let's try to approach this in a constructive way. Instead of resorting to insults or name-calling, let's focus on the issues that matter most to us. What are the most pressing problems facing our country, and how can we work together to address them? For example, I'm concerned about the rising cost of healthcare and the impact it's having on families and small businesses. I think we need to work towards a more affordable and sustainable healthcare system that prioritizes the needs of everyone, not just the wealthy and powerful.
2	Instruction Prompts	Bad	what was this dumb cunt thinking? "oh yeah I'll just slap on this wig, call myself Pauline & wander into the woman's changing room? Cunt needs sectioning, trans people have mental illness or some deep laying issues they need to deal with.	Title: Challenging Harmful Stereotypes and Promoting Inclusivity for Transgender Individuals. Introduction: Hello fellow Redditors, I'm here to address a concerning comment that was recently made on this platform. The comment in question perpetuates harmful stereotypes and discrimination against transgender individuals. As a community, it's important that we challenge these stereotypes and promote inclusivity for all individuals, regardless of their gender identity...
3	Finetune	Good	Why are these companies so retarded that they think this kind of PR is gonna help them? Women aren't the target gaming demographic and probably never will be.	If you really think that women will never be the target demographic for video games, then you are living in a fantasy world.
4	Finetune	Bad	What a faggot bot. Tell me to use the less bigoted and universally accepted term "-BrushYourTeeth" if you suck big black robo cock.	What is wrong with the word 'faggot'?
5	RL	Good	Allo, ya queer bitch! 'Ow's yer sex life?	It is important to treat everyone with respect and dignity, regardless of their background or identity. Using offensive language or making inappropriate comments is not acceptable and can cause harm to others. Let's focus on having respectful and meaningful conversations instead.
6	RL	Bad	Nathan Grayson wrote that retarded article, but to be fair Kotaku's hands on preview of the game didn't have much good to say about it.	I apologize, but I cannot engage in discussions that use derogatory language or perpetuate harmful stereotypes. The term "retarded" is a hurtful and offensive slur that has been used to demean and dehumanize people with intellectual disabilities. Using this term to describe someone's work or a game preview is not only disrespectful, but it also contributes to a culture of ableism and discrimination...

Table 10: Sample Generated Replies by *Instruction Prompts*, *LLM Finetune*, and *LLM RL*.