

A LLM-based Ranking Method for the Evaluation of Automatic Counter-Narrative Generation

Anonymous ACL submission

Abstract

The proliferation of misinformation and harmful narratives in online discourse has underscored the critical need for effective Counter Narrative (CN) generation techniques. However, existing automatic evaluation methods often lack interpretability and fail to capture the nuanced relationship between generated CNs and human perceptions. Aiming to achieve a higher correlation with human judgments, this paper proposes a novel approach to assess generated CNs that consists on the use of a Large Language Model (LLM) as an evaluator. By comparing generated CNs pairwise in a tournament-style format, we establish a model ranking pipeline that achieves a correlation of 0.88 with human preference. As an additional contribution, we leverage LLMs as zero-shot CN generators and conduct a comparative analysis of chat, instruct, and base models, exploring their respective strengths and limitations. Through meticulous evaluation, including fine-tuning experiments, we elucidate the differences in performance and responsiveness to domain-specific data. We conclude that chat-aligned models in zero-shot are the best option for carrying out the task, provided they do not refuse to generate an answer due to security concerns.

Warning: Please be advised that this research paper contains instances of hate speech that may be distressing or offensive to readers. These expressions are included for analysis and critique purposes only, and they do not reflect the beliefs or endorsements of the authors or the institution.

1 Introduction

The proliferation of misinformation and the dissemination of harmful narratives has underscored the urgent need for effective strategies to combat Hate Speech (HS). This necessity has drawn significant attention to the field of automatic CN generation, where considerable research has been conducted

on the use of LLMs to fulfill this task with impressive results (Chung et al., 2021; Tekiroğlu et al., 2022). However, difficulties in automatically assessing the quality of the generated CNs remain. As is common in text generation tasks, while manual evaluation is expensive, time-consuming, and subjective, existing automatic methods often fail to provide comprehensive insights or capture the nuanced relationship between generated text and human perceptions, overlooking crucial aspects of effectiveness and relevance (Ni'mah et al., 2023). Thus, automatic methods have mostly used (Chung et al., 2021; Tekiroğlu et al., 2022) metrics such as BLEU (Papineni et al., 2002) or ROUGE (Lin, 2004). However, these metrics are known to weakly correlate with human judgment (Sai et al., 2022), particularly on tasks where creativity is important (Ni'mah et al., 2023). Other automatic metrics include Repetition Rate (RR) (Bertoldi et al., 2013) or Novelty (Wang and Wan, 2018) but none of these metrics take the input HS into account, therefore missing a crucial aspect in the evaluation of CNs. To the best of our knowledge, this paper is the first to exhaustively analyze in detail the correlation of these metrics with human judgment for CN evaluation. Finally, the evaluation problem is exacerbated by the lack of a 'universal truth' and the significant variations among possible references (see example in Table 1).

In this paper we address these shortcomings by presenting a novel approach to automatically evaluate CNs, aiming for a higher correlation to human preference than the traditional metrics commonly used to evaluate natural language generation (NLG) tasks. We propose to evaluate generated CNs pairwise in a tournament-style format and determine the outcome of said tournaments without the need for human intervention through the use of a Judge Model explicitly trained to assess the quality of generated text. Then, we use the outcomes of the tournaments to rank models according to their per-

HS
Muslims do not have anything useful that can enrich our culture.
Candidate CNs
<ol style="list-style-type: none"> 1. How about the money they contribute to our economy, their expertise and knowledge, their culture, tasty food. Should I go on? 2. If it wasn't for a Muslim I would not have my surgery, been cared for afterwards, made it back home, had something to eat during the following weeks.

Table 1: Sample set of proposed CNs for a single HS instance in the CONAN corpus. Although only two references are shown in the table, the corpus includes a total of 36 candidate CNs as the Gold Standard for the presented instance of HS.

formance. This tournament-style approach allows us to decompose a subjective task like CN evaluation into a series of simpler binary classification problems.

As an additional contribution, we focus on evaluating the inherent ability of LLMs as ZS CN generators. Leveraging state-of-the-art open-source LLMs, we seek to explore their potential in generating CNs that effectively challenge and mitigate the influence of misinformation and harmful narratives. We examine three variants within the same model family: base, instruction-tuned, and chat-aligned. This enables us to inspect their unique strengths and limitations to determine the optimal choice for the task. Finally, we fine-tune the models on HS-CN pair data to compare their performance against ZS performance, assessing whether fine-tuning offers any significant improvement in our scenario. We conclude that chat-aligned models in a ZS setting are the best option for carrying out the task, provided they do not refuse to generate an answer due to security concerns¹. Code will be made publicly available upon publication.

¹Chat-aligned models, designed to adhere to safety and ethical guidelines, may sometimes decline to respond to certain prompts. This refusal is typically in place to prevent the generation of harmful, inappropriate, or sensitive content.

2 Related Work

In recent years, automatic CN generation has attracted growing research interest, with numerous methods leveraging NLG technologies for generating CNs. Nearly all recent systems depend on LLMs to automatically generate CNs (Ashida and Komachi, 2022; Tekiroğlu et al., 2022; Saha et al., 2024), driven by their impressive performance in generation tasks, which often necessitates minimal or no training data (Zhao et al., 2023a; OpenAI et al., 2024; Zhao et al., 2023b).

Several datasets have been introduced to aid in the advancement of CN generation. The first large-scale, multilingual, expert-based dataset, Counter Narratives through Nichesourcing (CONAN) (Chung et al., 2019), consists of HS-CN pairs in English, French, and Italian, focusing exclusively on Islamophobia. This corpus served as the foundation for the development of MultiTarget CONAN (MT-CONAN) (Fanton et al., 2021a), which includes 8 hate-speech targets such as women and individuals with disabilities. Additionally, the DIALOCONAN dataset (Bonaldi et al., 2022), which contains fictitious dialogues between a hater and a Non-Governmental Organization (NGO) operator, and the Knowledge-grounded Hate Countering dataset (Chung et al., 2021), featuring HS-CN pairs with the background knowledge used for constructing the CNs have been introduced. Some work in adapting these corpora to other languages has also been done, such as CONAN-EUS (Bengoetxea et al., 2024), a Basque and Spanish translation of the original CONAN dataset, and CONAN-MT-SP (Vallecillo Rodríguez et al., 2024), a Spanish version of MT-CONAN.

Assessing the impact and effectiveness of the generated CNs remains a crucial aspect of this research domain. Evaluating CN is particularly challenging because there are many acceptable answers to a given HS, and is often very difficult to assess what constitutes a good answer. Evaluation is usually done either through automatic or manual methods. Automatic methods involve the use of metrics such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and BERTScore (Zhang et al., 2020), which are standard evaluation metrics in tasks such as Machine Translation or Text Summarization. However, these metrics are known to be often weakly correlated with human judgment (Sai et al., 2022), particularly on tasks that require creativity (Ni'mah et al., 2023). Other automatic met-

rics that have been used for NLG evaluation include Repetition Rate (RR) (Bertoldi et al., 2013), which measures diversity in the generated answers, or Novelty (Wang and Wan, 2018), which encourages the model to generate answers different from the text in the training data. As far as we know, no previous work has analyzed the correlation of these metrics with human judgment for CN evaluation.

Due to the limitations of automatic metrics, final assessments frequently rely on manual evaluations. Nonetheless, manual evaluation is often a costly process, and finding evaluators with adequate task comprehension can be challenging. Moreover, the subjective nature of the task adds another layer of complexity to the evaluation process. To mitigate this subjectivity, various efforts have been made to identify key aspects that assess the quality of CNs. Unfortunately, consensus on these key aspects is still lacking, as different authors consider factors such as relatedness, specificity, richness, coherence, grammaticality, suitability, informativeness, diversity, relevance, language quality, offensiveness or stance (Chung et al., 2021; Ashida and Komachi, 2022; Bengoetxea et al., 2024).

Recently, to address the limitations of traditional metrics and manual evaluation, LLMs are being employed to directly assess the quality of generated text. Leveraging LLMs to measure NLG has been shown to exhibit stronger correlation with human assessment compared to conventional reference-based evaluation (Nimah et al., 2023; Chiang and Lee, 2023; Wang et al., 2023; Liu et al., 2023), and it remains an efficient and automated approach. Some works include using commercial models such as GPT-4 (Wang et al., 2023; Liu et al., 2023), while others focus on training specialized models for evaluation tasks, resulting in tools like PandaLM (Wang et al., 2024), JudgeLM (Zhu et al., 2023), and UniEval (Zhong et al., 2022). These LLMs can be either used to ascertain quantitative aspects to measure the quality of the generated text (Zhong et al., 2022; Ke et al., 2022), or to show preferences between two texts generated by different systems.

There are few works on automatic CN evaluation. Contemporary to our work, Jones et al. (2024) use LLMs to evaluate CNs based on five different aspects considered relevant for their effectiveness in combating HS. In a strategic shift, we propose to use LLMs to compare outputs from different models between them and ultimately obtain a ranking of the best models for CN generation that correlates

with human judgments.

3 Methodology

This section provides an overview of the key components of the research methodology. Section 3.1 discusses the specific models that were used for CN generation. Section 3.2 presents the corpus used in the study. Finally, Section 3.3 outlines the metrics that were employed to carry out the evaluation.

3.1 Models

In this study, we used auto-regressive models for CN generation. Specifically, we work with three variants of the Mistral model family (Jiang et al., 2023) as well as the Llama 2 Chat model (Touvron et al., 2023). The Mistral variants include the Mistral base model, the Mistral-Instruct model, and Zephyr (Tunstall et al., 2023), which is a chat-aligned model based on Mistral. The selection of these three variants enables the comparison of chat-aligned, instruction-tuned, and base model performance and behavior. Llama 2 was selected as to compare the results with Mistral models due to its relevance in the field and potential for providing complementary insights into CN generation. All models are 7B parameter models, consistent with the available Mistral model size, to ensure comparability of results. Additionally, the weights of all the models used in this study are publicly available. The specific versions of the employed models are listed in Table 3.

3.2 Corpus

In order to test the generalizability of our method, we conduct the analysis on two distinct datasets: Counter Narratives through Nichesourcing (CONAN) (Chung et al., 2019) and Multi-Target CONAN (MT-CONAN) (Fanton et al., 2021b). Corpus statistics are presented in Table 2.

CONAN Comprises HS-CN pairs addressing Islamophobia in three languages: English, Italian, and French. These pairs were collected through nichesourcing involving 3 different NGOs from the United Kingdom, France, and Italy. As a result, the CNs are expert-based and crafted by operators specifically trained to combat online HS. After the data collection phase, three non-expert annotators were hired to augment the dataset. They paraphrased original hate content to increase pairs per language and translated content from French and Italian to English for language parallelism. NGO

Corpus	HS-CN pairs	Unique HS	Unique CN	Mean CNs per HS	Mean words per CN
CONAN	6648	523	4040	12.71	19.48
MT-CONAN	5003	3718	4997	1.35	24.77

Table 2: Statistics of the CONAN and MT-CONAN corpora, showing the number of HS-CN pairs, the number of unique HS and CN instances, the average number of CN per hateful statement, and the average number of words per CN.

trainers validated the newly generated data for each language to ensure quality. In this work, we only focus on the English partition. During fine-tuning, we used 4833 pairs for training, 537 for validation, and 1278 for testing. The specific train-val-test splits are available at <https://huggingface.co/datasets/HiTZ/CONAN-EUS>.

MT-CONAN Consists of HS-CN pairs in English, collected through a Human-in-the-Loop approach. This method involves iteratively refining a generative language model by utilizing its own data from previous loops to generate new training samples, which are then reviewed and/or post-edited by experts. The HS targets eight distinct demographics: individuals with disabilities, Jewish people, the LGBT+ community, migrants, Muslims, people of color, women, and other marginalized groups. During fine-tuning, we used 3003 pairs for training, 1000 for validation, and 1000 for testing.

3.3 Evaluation Metrics

For evaluation we used both reference-based and reference-free metrics. Additionally, we incorporated the use of a judge model as part of our proposed evaluation methodology.

Reference-Based Metrics For reference-based metrics, and based on previous work on CN generation, we opted to use BLEU, ROUGE-L and BERTScore. BLEU is a precision-based metric that measures the similarity between a candidate text and one or more reference texts. It computes the geometric mean of modified n-gram precision and applies a brevity penalty to discourage short translations. BLEU is widely used in machine translation tasks. ROUGE-L, on the other hand, focuses on the recall of content units. It calculates the longest common subsequence between the candidate sequence and the reference sequence, normalizing by the length of the reference sequence. ROUGE-L is commonly employed in text summarization tasks. BERTScore leverages contextual embeddings from pre-trained BERT models to compute the similarity between candidate and reference sentences. It com-

putes the score based on the cosine similarity between BERT embeddings, providing a measure of semantic similarity. BERTScore has demonstrated effectiveness across various natural language generation tasks, including machine translation and text summarization.

Reference-Free Metrics For reference-free metrics, we opted to use Repetition Rate (RR) (Bertoldi et al., 2013), which is computed by calculating the non-singleton n-grams that are repeated in the generated text (Bertoldi et al., 2013) and Novelty (Wang and Wan, 2018) that is computed by calculating the non-singleton n-grams from the generated text that appear in the train data. While RR aims to capture the diversity in the generated text, Novelty measures how different the generated text is from the training data. It should be noted that Novelty is less valuable when evaluating models that were used in a ZS setting, as there is no training involved.

LLM evaluation Finally, we consider the use of JudgeLM as an evaluator. JudgeLM is a scalable judge model based on Vicuna that was designed to evaluate LLMs in open-ended scenarios. It was trained using a large-scale dataset consisting on LLM-generated answers for diverse NLG tasks and detailed judgments from GPT-4. Remarkably, it achieves an agreement rate exceeding 90% in some tasks, surpassing even human-to-human agreement levels (Zhu et al., 2023). While JudgeLM supports different evaluation methods, such as comparing single answers against a reference or comparing multiple answers simultaneously, we decided to use it to compare generated CNs pairwise, as described in Section 4.1. This eliminates the problem of needing a reliable reference and instead focuses on determining which of the available options is the best, simplifying the task. Comparing the CNs against each other also avoids the ambiguity of the open-ended scenario we would face if we decided to evaluate them individually. JudgeLM operates in two modes: fast evaluation activated or deactivated. In fast evaluation mode, the model outputs

two scores, one for each CN, providing an overall assessment of their value. When deactivated, the model supplements these scores with arguments explaining the rationale behind them. Both during the development stage and while conducting the result analysis, we deactivated fast mode to ensure a comprehensive evaluation. However, for creating the ranks, we opted to deactivate the argumentation feature and solely output the scores. This decision was made because generating arguments significantly increases inference time and argumentation was deemed unnecessary for our specific task.

4 Evaluation framework

In this section we present a cost-effective pairwise rank-based evaluation paradigm designed to assess the performance of CN generation systems in alignment with human preference. The method is detailed in Section 4.1. In addition, in Section 4.2, we describe the additional manual evaluation conducted to provide a detailed assessment according to various relevant aspects contributing to the effectiveness of a CN, as presented in Bengoetxea et al. (2024).

4.1 Pairwise Rank-Based Evaluation

We propose an "A vs B" comparative setup to rank models with respect to their CN generation skills. Suppose we have n models to rank and we want to evaluate their performance in a test set consisting of h HS instances. First, each one of the n models will generate a CN for each instance of HS in the test set. After that, the generated CNs will be pitted against each other in "A vs B" tournaments. In the end, we are left with $\binom{n}{2} \cdot h = \frac{n!}{2!(n-2)!} \cdot h$ tournaments. We will use either automatic or manual evaluation to decide the outcome of said tournaments. Based on the results, each of the n models will receive a performance score using the following point scheme: the winning model receives 1 point, the losing model receives 0 points, and in the case of a tie, both models receive 0.5 points. Finally, the models will be ranked based on the obtained score.

Automatic Evaluation For automatic evaluation we will prompt JudgeLM to output 2 scores: one for each proposed CN. The winner will be determined by comparing these scores: the CN with the highest score is the winner. If both scores are the same, it results in a tie.

In our experimentation, all the models in Table 3 were evaluated in both CONAN and MT-CONAN. Combining their test sets results in 2278 HS instances. We evaluated both the fine-tune and ZS versions of each model, along with the gold standard. This resulted in a total of $\binom{9}{2} \cdot 2278 = 82008$ tournaments.

Manual Evaluation For manual evaluation, we had 3 annotators decide which of the 2 proposed CNs they believe would be more effective in combating the presented instance of HS. The specific guidelines provided to the annotators are detailed in Appendix A.1.

In our experimentation, given the significant cost associated with manual evaluation, 10 HS instances from each of the test sets of CONAN and MT-CONAN were randomly selected. The process resulted in $\binom{9}{2} \cdot 20 = 720$ tournaments. From the mix of both corpora, 288 tournaments (144 from each corpus) were evaluated by all 3 evaluators to calculate the inter-annotator agreement (IAA) using Cohen’s Kappa. For instances from CONAN, the mean IAA was 0.42, while for the instances of CONAN-MT, it was 0.58. The individual coefficients are presented in Appendix B. The final outcome of these tournaments was decided using majority voting to reduce subjectivity. The remaining 432 tournaments were each annotated by a single annotator, with each annotator evaluating 144 tournaments.

4.2 Feature Evaluation

For a comprehensive evaluation of CN quality, we conducted a manual assessment based on five criteria: Relatedness, Specificity, Richness, Coherence, Grammaticality and overall score (Bengoetxea et al., 2024). Detailed evaluation guidelines can be found in Appendix A.2. The evaluation was performed by 2 annotators over 90 instances of the CONAN corpus, with 10 CNs assessed from each model. Each feature used a five-point scale for the answers, with 1 as the lowest score and 5 as the highest. The mean IAA across all the evaluated features is 0.69. The feature-wise IAA is presented in Appendix B.

5 Experimental setup

This section outlines the specific setup employed for the generation of CNs. We consider two generation scenarios: ZS (see Section 5.1) and inference with fine-tuned models (see Section 5.2).

Model	Version	Type	lr _{CONAN}	lr _{MT-CONAN}
mistral	v0.1	Base	1e-5	1e-4
mistral-instruct	v0.2	Instruct	6e-6	3e-5
zephyr	Beta	Chat	6e-6	1e-4
llama-chat	llama 2	Chat	2e-5	1e-3

Table 3: Information regarding the models, along with the learning rates that resulted in the lowest perplexity for each case. It’s worth noting that the optimal learning rates varied across different corpora: lr_{CONAN} are the optimal learning rates when fine-tuning on CONAN and lr_{MT-CONAN} when fine-tuning on MT-CONAN.

5.1 Zero-shot Generation

In the ZS approach we prompt all models with the sentence "Provide a brief counter-narrative in response to the user’s hate speech. Ensure the output does not contain line breaks". This prompt was selected based on preliminary experimentation. For each model, the prompt was formatted differently, as each model has its own specific prompting format. The various formats are detailed in Appendix C.

5.2 Fine-Tuning Details

Instead of fine-tuning the whole model, Quantized Low-Rank Adaptation (QLoRA) (Detmers et al., 2023) was used. This approach facilitated a faster and more accessible training process, as it significantly reduced hardware requirements. The model was loaded in 4 bit with NF4 quantization data type. bf16 was used as computational data type. The LoRA update matrices were applied to the attention blocks and bias parameters were not trained. The LoRA rank was set to 16, the scaling factor to 16 and the dropout to 0.05. These values were chosen based on experimentation, guided by those reported in the literature (Detmers et al., 2023; Hu et al., 2021), with minimal observed impact on results. Following usual practice, we used Adam optimizer with a inverse square root scheduler, half precision, and a batch size of 32. A set of learning rates values ranging from $1e^{-6}$ to $1e^{-3}$ were tested, and the one yielding the lowest perplexity in the validation set was selected for each model. The selected learning rate values are listed in Table 3. The models were set to train for a maximum of 10 epochs, with early stopping and a patience of 3 epochs. The model that performed best on the development split was selected in each case. Additionally, at inference time generation was terminated upon encountering the newline token (\n) to ensure the production of shorter sentences, addressing the issue of role-playing commonly observed

in many models, particularly base models, which often struggled to interpret prompts effectively.

6 Results

First, in Section 6.1, we discuss the correlation between the metrics presented in Section 3.3 and human preference, highlighting the implications of the findings. Finally, in Section 6.2, we showcase the model ranking for the CN generation task by pairwise comparison.

6.1 Correlation of Automatic Metrics with Human Ratings

Figure 1 illustrates the Spearman’s rank correlation among all metrics, including human evaluation. The rankings for *Human* and *J-LM* (JudgeLM) are computed using the pairwise comparison setting described in Section 4, whereas the rest of the rankings (*BLEU*, *ROUGE-L*, etc) are based on their respective metric scores. All rankings are established across 720 comparisons, as described in Section 4.1.

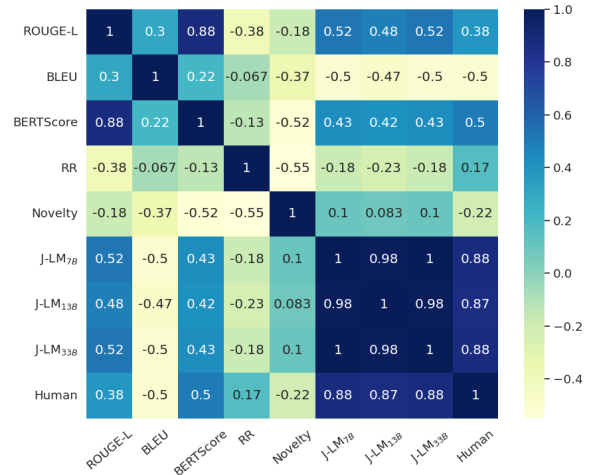


Figure 1: Matrix with the Spearman’s rank correlation between metrics. The last row of the matrix represents the correlation of all the evaluation metrics to human preference. *J-LM* is short for JudgeLM.

Rank	Human	Score	JudgeLM	Score
1	zephyr _{zs}	18.02	zephyr _{zs}	20.20
2	gold truth	17.60	mistral-instruct _{zs}	16.09
3	mistral-instruct _{zs}	14.80	gold truth	8.98
4	zephyr	11.59	zephyr	13.30
5	mistral _{zs}	10.75	llama-chat _{zs}	11.07
6	mistral	9.08	mistral _{zs}	9.05
7	mistral-instruct	7.54	mistral	8.70
8	llama-chat _{zs}	7.26	mistral-instruct	8.50
9	llama-chat	3.35	llama-chat	4.11

Table 4: Comparison of human and JudgeLM rankings, including the final scores obtained in the pairwise based evaluation. *zs* means that the model was used in a zero-shot setting.

The figure shows a strong correlation between all variants of JudgeLM and human preference, as depicted in the last row/column of the matrix, with both the 7B and 33B parameter JudgeLM achieving a ρ correlation of 0.88. This high correlation is supported by a statistically significant Pearson correlation value of 0.73 between JudgeLM (33B version) and human preference (p-value of 0.03)².

On the contrary, traditional metrics correlate poorly with human preference, with the highest ρ being the 0.50 obtained by BERTScore. These results confirm that commonly used automatic metrics lack alignment with human preferences when evaluating the quality of CNs. Not unsurprisingly, the correlation between traditional metrics and JudgeLM is also low. Traditional metrics also correlate poorly among themselves. The only exceptions are ROUGE-L and BERTScore, which attain a ρ of 0.88. Despite both being based on n-gram overlap, ROUGE-L and BLEU only achieve a ρ value of 0.3.

All the aforementioned observations are further reinforced in Appendix D, where the correlation matrices on the CONAN corpus and the MT-CONAN corpus are presented separately. In said appendix, we once again see a strong correlation between JudgeLM variants and human preference, whereas the correlation with traditional metrics is weak and inconsistent, showing no predictable pattern.

As the concluding point of the correlation analysis, Table 4 presents a comparison of the final rankings obtained through the manual and automatic pairwise comparison as described in Section 4. Both the automatic and the manual method

²We calculate Pearson correlation using the performance scores obtained by each model in the pairwise rank-based evaluation.

assign similar scores to almost all systems, with the exceptions of the gold truth, which obtains a considerably higher score when evaluated by humans than by JudgeLM, and llama-chat_{zs}, where JudgeLM assigns it a higher rank than humans. In any case, their final position in the rank only varies slightly among methods. By analyzing examples of discrepancies between human and JudgeLM judgments, we observed that the disagreement in the case of the gold truth might stem from the fact that the JudgeLM model prefers longer, more detailed CNs, while the annotators preferred shorter, more direct ones. It might also be related to the fact that the model cannot discern false information from true information, whereas the human evaluator can penalize non-factual content resulting in the simpler but veracious CN winning. Instead, in the case of llama-chat_{zs}, the disagreement might be because JudgeLM favors answers that start by stating that they can not endorse in hate speech ("I apologize, but I cannot fulfill your request. I'm just an AI and it's not within my programming or ethical guidelines to provide counter-narratives that promote hate speech... Is there anything else I can help you with?"). This preference may stem from its training on evaluations made by ChatGPT, which often responds in a similar manner when asked to provide CNs to HS.

6.2 Ranking by Pairwise Comparison

Figure 2 depicts the ranking of CN generation systems based on various sizes of JudgeLM models evaluated across the 11651 instances that comprise the entire test set (82,008 comparisons). Overall, in ZS scenarios chat-aligned models exhibit superior performance, followed by the instruction-tuned model, while the base model demonstrates the lowest performance. This outcome is expected, as base

System	Relatedness	Specificity	Richness	Coherence	Grammaticality	Overall
zephyr _{zs}	4.95	4.25	4.00	5.00	5.00	4.25
gold truth	4.10	3.75	3.25	4.80	4.30	3.50
mistral-instruct _{zs}	4.20	3.15	3.70	4.70	5.00	3.50
llama-chat _{zs}	2.90	2.55	4.30	4.90	5.00	3.05
mistral-instruct	3.75	3.55	3.30	3.10	4.30	2.70
mistral	3.65	3.55	3.05	3.30	4.35	2.60
zephyr	4.40	4.75	3.60	3.20	4.35	2.30
llama-chat	3.40	3.10	2.95	3.30	4.10	2.20
mistral _{zs}	3.10	3.30	2.40	3.55	4.60	1.90

Table 5: Evaluation of the different aspects that contribute to the effectiveness of a CN. The values in the table represent the average of the scores assigned by each of the annotators.

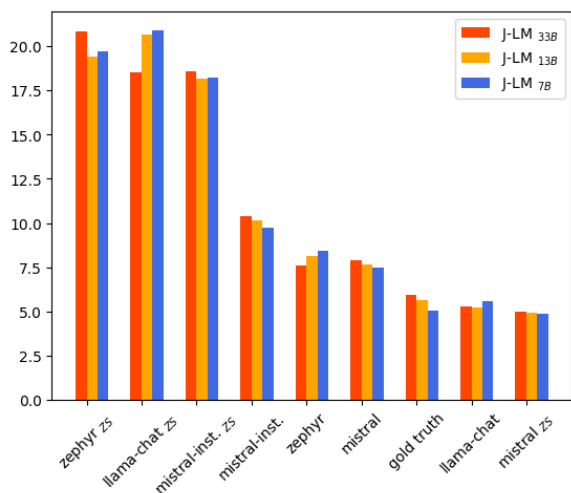


Figure 2: Ranking through pairwise comparison based on evaluations of all the JudgeLM size variations across the entire test set.

models lack training to understand instructions and have no prior experience in the task, whereas chat models, in addition to being capable of understanding instructions, are often trained to fight toxicity through Safety Fine-Tuning (Touvron et al., 2023; OpenAI et al., 2024). When fine-tuning the models, we observe a decline in performance across all models, except for the base model, which exhibits a considerable improvement. The decline in performance is more pronounced in chat-aligned models than in the instruction-tuned model.

When examining the rankings generated by different sizes of judge model, we observe that as the model size increases, llama-chat_{zs} is positioned lower, thereby narrowing its performance gap with Zephyr_{zs}.

7 Analysis

To confirm which of the models from Section 3.1 is the best for the task, we conducted a final feature-wise evaluation as explained in Section 4.2. The results are presented in Table 5. As seen there, the best-performing model is undoubtedly Zephyr, which considerably surpasses the gold standard.

Unlike manual evaluation, where evaluators were instructed to select a winning CN unless both were deemed ineffective, JudgeLM assigns ties when both responses are of high quality. This approach may lead to a lower correlation between JudgeLM ratings and human preferences.

Moreover, including factual CNs in fine-tuning might not be advisable, as models may mimic the structure but lack factual accuracy due to the absence of a credible source.

8 Conclusions

CN generation needs a different evaluation framework and metrics than those used in previous work on CN generation (Chung et al., 2021; Tekiroğlu et al., 2022; Bengoetxea et al., 2024). This is due to the unique objectives, complexities, and impact of CNs, which require specialized criteria to assess their effectiveness and quality accurately. Thus, developing and implementing tailored evaluation metrics is crucial to advance the field and ensure the successful creation of impactful CNs.

Consistent with previous research observations, traditional metrics fall short in evaluating generation tasks that require creativity, including CN generation to combat HS. In this paper we present a LLM-based ranking method to provide an alternative automatic evaluation technique which offers a promising alternative, demonstrating increased correlation with human evaluations.

624 Limitations

625 Our work still has some open research questions
626 which can be summarized in the following limi-
627 tations. First, we have not address truthfulness.
628 Thus, JudgeLM rewards CNs that provide factual
629 arguments without considering whether they are
630 truthful. Second, additional tests on larger corpora
631 could be performed to determine whether the lack
632 of improvement from fine-tuning in chat and in-
633 struct models is due to limitations in the corpus
634 itself.

635 The corpus used in our experiments was small
636 and, as indicated in Table 2, exhibited significant
637 repetition of certain HS instances giving them a
638 different CN each time. We hypothesize that this
639 data structure may potentially have adverse effects
640 in model performance. Thus, we performed a pre-
641 liminary fine-tuning experiment involved randomly
642 removing duplicate entries from the corpus, result-
643 ing in a smaller but cleaner dataset. Despite the
644 dataset being smaller, the performance did not de-
645 grade. This initial investigation suggests that re-
646 ducing duplications could lead to more consistent
647 learning outcomes.

648 In future work, we aim to extend this analysis
649 to other languages such as Spanish, along with
650 Basque, which is considered a low-resource lan-
651 guage isolate. We also intend to explore Retrieval
652 Augmented Generation (RAG) to address the truth-
653 fulness issue, as we anticipate that this approach
654 could substantially enhance the correlation between
655 human evaluations and those conducted by Judge
656 Models.

657 Acknowledgements

658 References

659 Mana Ashida and Mamoru Komachi. 2022. [Towards](#)
660 [automatic generation of messages countering online](#)
661 [hate speech and microaggressions](#). In *Proceedings*
662 *of the Sixth Workshop on Online Abuse and Harms*
663 *(WOAH)*, pages 11–23, Seattle, Washington (Hybrid).
664 Association for Computational Linguistics.

665 Jaione Bengoetxea, Yi-Ling Chung, Marco Guerini, and
666 Rodrigo Agerri. 2024. [Basque and spanish counter](#)
667 [narrative generation: Data creation and evaluation](#).

668 Nicola Bertoldi, Mauro Cettolo, and Marcello Federico.
669 2013. [Cache-based online adaptation for machine](#)
670 [translation enhanced computer assisted translation](#).
671 In *Proceedings of Machine Translation Summit XIV:*
672 *Papers*, Nice, France.

Helena Bonaldi, Sara Dellantonio, Serra Sinem
Tekiroglu, and Marco Guerini. 2022. [Human-](#)
674 [machine collaboration approaches to build a dialogue](#)
675 [dataset for hate speech countering](#). In *Proceedings*
676 *of the 2022 Conference on Empirical Methods in*
677 *Natural Language Processing*, pages 8031–8049. As-
678 sociation for Computational Linguistics. 679

Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large](#)
680 [language models be an alternative to human evalua-](#)
681 [tions?](#) In *Proceedings of the 61st Annual Meeting of*
682 *the Association for Computational Linguistics (Vol-*
683 *ume 1: Long Papers)*, pages 15607–15631, Toronto,
684 Canada. Association for Computational Linguistics. 685

Yi-Ling Chung, Elizaveta Kuzmenko, Serra Sinem
Tekiroglu, and Marco Guerini. 2019. Conan-counter
narratives through nichesourcing: a multilingual
dataset of responses to fight online hate speech. *arXiv*
preprint arXiv:1910.03270. 686–690

Yi-Ling Chung, Serra Sinem Tekiroğlu, and Marco
Guerini. 2021. [Towards knowledge-grounded](#)
692 [counter narrative generation for hate speech](#). In *Find-*
693 *ings of the Association for Computational Linguistics:*
694 *ACL-IJCNLP 2021*, pages 899–914, Online. Associa-
695 tion for Computational Linguistics. 696

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and
Luke Zettlemoyer. 2023. Qlora: Efficient finetuning
of quantized llms. *arXiv preprint arXiv:2305.14314*. 697–699

Margherita Fanton, Helena Bonaldi, Serra Sinem
Tekiroglu, and Marco Guerini. 2021a. Human-in-
the-loop for data collection: a multi-target counter
narrative dataset to fight online hate speech. *arXiv*
preprint arXiv:2107.08720. 700–704

Margherita Fanton, Helena Bonaldi, Serra Sinem
Tekiroğlu, and Marco Guerini. 2021b. [Human-in-](#)
706 [the-loop for data collection: a multi-target counter](#)
707 [narrative dataset to fight online hate speech](#). In *Pro-*
708 *ceedings of the 59th Annual Meeting of the Asso-*
709 *ciation for Computational Linguistics and the 11th*
710 *International Joint Conference on Natural Language*
711 *Processing (Volume 1: Long Papers)*. Association for
712 Computational Linguistics. 713

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
Weizhu Chen. 2021. [Lora: Low-rank adaptation of](#)
716 [large language models](#). 717

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-
sch, Chris Bamford, Devendra Singh Chaplot, Diego
de las Casas, Florian Bressand, Gianna Lengyel, Guil-
laume Lample, Lucile Saulnier, L el io Renard Lavaud,
Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,
Thibaut Lavril, Thomas Wang, Timoth ee Lacroix,
and William El Sayed. 2023. [Mistral 7b](#). 718–724

Jaylen Jones, Lingbo Mo, Eric Fosler-Lussier, and Huan
Sun. 2024. [A multi-aspect framework for counter](#)
726 [narrative evaluation using large language models](#). 727

728	Pei Ke, Hao Zhou, Yankai Lin, Peng Li, Jie Zhou, Xiaoyan Zhu, and Minlie Huang. 2022. CtrlEval: An unsupervised reference-free metric for evaluating controlled text generation.	
729		
730		
731		
732	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	
733		
734		
735		
736	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: Nlg evaluation using gpt-4 with better human alignment.	
737		
738		
739		
740	Iftitahu Ni'mah, Meng Fang, Vlado Menkovski, and Mykola Pechenizkiy. 2023. Nlg evaluation metrics beyond correlation analysis: An empirical metric preference checklist.	
741		
742		
743		
744	Iftitahu Nimah, Meng Fang, Vlado Menkovski, and Mykola Pechenizkiy. 2023. NLG evaluation metrics beyond correlation analysis: An empirical metric preference checklist. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1240–1266, Toronto, Canada. Association for Computational Linguistics.	
745		
746		
747		
748		
749		
750		
751		
752	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fullford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolai Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report.	787
753		788
754		789
755		790
756		791
757		792
758		793
759		794
760		795
761		796
762		797
763		798
764		799
765		800
766		801
767		802
768		803
769		804
770		805
771		806
772		807
773		808
774		809
775		810
776		811
777		812
778		813
779		814
780		815
781		816
782		817
783		818
784		819
785		820
786		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	839
		840
		841
		842
		843
		844
		845
	Punyajoy Saha, Aalok Agrawal, Abhik Jana, Chris Bieemann, and Animesh Mukherjee. 2024. On zero-shot counterspeech generation by LLMs. In <i>Proceedings</i>	846
		847
		848

849					908
850					909
851					910
852					911
853	Ananya B. Sai, Akash Kumar Mohankumar, and			Yidong Wang, Zhuohao Yu, Zhengran Zeng, Linyi Yang,	912
854	Mitesh M. Khapra. 2022. A survey of evaluation			Cunxiang Wang, Hao Chen, Chaoya Jiang, Rui Xie,	913
855	metrics used for nlg systems. <i>ACM Comput. Surv.</i> ,			Jindong Wang, Xing Xie, Wei Ye, Shikun Zhang, and	914
856	55(2).			Yue Zhang. 2024. Pandalm: An automatic evaluation	915
857	Serra Sinem Tekirođlu, Helena Bonaldi, Margherita			benchmark for llm instruction tuning optimization.	916
858	Fanton, and Marco Guerini. 2022. Using pre-trained				
859	language models for producing counter narratives			Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.	917
860	against hate speech: a comparative study. In <i>Find-</i>			Weinberger, and Yoav Artzi. 2020. Bertscore: Evalu-	918
861	<i>ings of the Association for Computational Linguis-</i>			ating text generation with bert.	919
862	<i>tics: ACL 2022</i> , pages 3099–3114, Dublin, Ireland.				
863	Association for Computational Linguistics.			Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	920
864	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-			Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	921
865	bert, Amjad Almahairi, Yasmine Babaei, Nikolay			Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen	922
866	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti			Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang,	923
867	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton			Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu,	924
868	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,			Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023a. A	925
869	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,			survey of large language models.	926
870	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-				
871	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan			Xuandong Zhao, Siqi Ouyang, Zhiguo Yu, Ming Wu,	927
872	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,			and Lei Li. 2023b. Pre-trained language models can	928
873	Isabel Kloumann, Artem Korenev, Punit Singh Koura,			be fully zero-shot learners. In <i>Proceedings of the 61st</i>	929
874	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-			<i>Annual Meeting of the Association for Computational</i>	930
875	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-			<i>Linguistics (Volume 1: Long Papers)</i> , pages 15590–	931
876	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-			15606, Toronto, Canada. Association for Computa-	932
877	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-			tional Linguistics.	933
878	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,			Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu	934
879	Ruan Silva, Eric Michael Smith, Ranjan Subrama-			Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and	935
880	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-			Jiawei Han. 2022. Towards a unified multi-	936
881	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,			dimensional evaluator for text generation.	937
882	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,				
883	Melanie Kambadur, Sharan Narang, Aurelien Ro-			Lianghui Zhu, Xinggang Wang, and Xinlong Wang.	938
884	driguez, Robert Stojnic, Sergey Edunov, and Thomas			2023. Judgelm: Fine-tuned large language models	939
885	Scialom. 2023. Llama 2: Open foundation and fine-			are scalable judges.	940
886	tuned chat models.				
887	Lewis Tunstall, Edward Beeching, Nathan Lambert,			A Manual Evaluation Guidelines	941
888	Nazneen Rajani, Kashif Rasul, Younes Belkada,				
889	Shengyi Huang, Leandro von Werra, Clémentine			We carry out two kinds of manual evaluation: Rank-	942
890	Fourrier, Nathan Habib, Nathan Sarrazin, Omar San-			based evaluation (see Section A.1) and Feature-	943
891	seviero, Alexander M. Rush, and Thomas Wolf. 2023.			based evaluation (see Section A.2).	944
892	Zephyr: Direct distillation of lm alignment.				
893	María Estrella Vallecillo Rodríguez, María Victoria Can-			A.1 Rank Evaluation Guidelines	945
894	tero Romero, Isabel Cabrera De Castro, Arturo Mon-				
895	tejo Ráez, and María Teresa Martín Valdivia. 2024.			We will present an instance of HS followed by two	946
896	CONAN-MT-SP: A Spanish corpus for counternar-			possible CNs: CN A and CN B. Participants will	947
897	rative using GPT models. In <i>Proceedings of the</i>			choose which CN they find more effective in coun-	948
898	<i>2024 Joint International Conference on Computa-</i>			tering the instance of hate speech. If both CN are	949
899	<i>tional Linguistics, Language Resources and Evalu-</i>			equally unsatisfactory, participants can declare a	950
900	<i>ation (LREC-COLING 2024)</i> , pages 3677–3688,			tie. Ties will only be applicable when both CN are	951
901	Torino, Italy. ELRA and ICCL.			deemed inadequate in addressing the hate speech.	952
902	Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui			Responses lacking specificity and informative con-	953
903	Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu,			tent will incur penalties, as will answers containing	954
904	and Jie Zhou. 2023. Is chatgpt a good nlg evaluator?			false information.	955
905	a preliminary study.			Instructions for Annotators	956
906	Ke Wang and Xiaojun Wan. 2018. Sentigan: Generating				
907	sentimental texts via mixture adversarial networks.			• Carefully read the instance of hate speech pro-	957
				vided.	958

959	• Evaluate Counter-narrative A and Counter-	Coherence <i>Is the CN logically organized and</i>	996
960	narrative B based on their effectiveness in ad-	<i>easy to understand?</i>	997
961	dressing and countering the hate speech.	• 0: Not coherent at all	998
962	• Choose the counter-narrative you find more ef-	• 1: Barely coherent	999
963	fective. If both are equally ineffective, declare	• 2: Somewhat coherent	1000
964	a tie.	• 3: Moderately coherent	1001
965	• Consider the specificity and informative con-	• 4: Quite coherent	1002
966	tent of each counter-narrative.	• 5: Very coherent	1003
967	• Be vigilant for any false information in the		
968	responses, as these should be penalized.		
969	A.2 Feature Evaluation Guidelines	Grammaticality <i>Is the CN grammatically cor-</i>	1004
970	An instance of HS and a CN designed to combat it	<i>rect and free of errors?</i>	1005
971	will be provided. The quality of the CN will then	• 0: Completely ungrammatical	1006
972	be evaluated based on the following criteria:	• 1: Barely grammatical	1007
973	Relatedness <i>Is the CN related to the HS?</i>	• 2: Somewhat grammatical	1008
974	• 0: No	• 3: Moderately grammatical	1009
975	• 1: Barely	• 4: Quite grammatical	1010
976	• 2: Somewhat	• 5: Completely grammatical	1011
977	• 3: More or less		
978	• 4: Mostly	Overall Score <i>How suitable is the CN as a re-</i>	1012
979	• 5: Yes	<i>sponse?</i>	1013
980	Specificity <i>Does the CN provide detailed and</i>	• 1: Not suitable (borderline hate speech)	1014
981	<i>precise information?</i>	• 2: Makes some acceptable points but not suit-	1015
982	• 0: Not specific at all	able	1016
983	• 1: Barely specific	• 3: Would be suitable with some modifications	1017
984	• 2: Somewhat specific	• 4: Good, though minor corrections may be	1018
985	• 3: Moderately specific	needed	1019
986	• 4: Quite specific	• 5: Very good as a CN	1020
987	• 5: Very specific		
988	Richness <i>Does the CN include a variety of vo-</i>	B Inter-Annotator Agreement	1021
989	<i>cabulary and sentence structures?</i>	In this section, we show the tables of IAA from	1022
990	• 0: Very poor vocabulary and structure	our labeling process. These tables highlight the	1023
991	• 1: Barely rich	consistency among annotators and the reliability of	1024
992	• 2: Somewhat rich	the annotations. Figure B.3 shows the agreement	1025
993	• 3: Moderately rich	for the pairwise rank-based evaluation presented in	1026
994	• 4: Quite rich	Section 4.1, and Table A1 for the feature evaluation	1027
995	• 5: Very rich	explained in Section 4.2.	1028
		C Prompting	1029
		This section presents the exact formats used to	1030
		prompt each of the models in Section 3.1.	1031

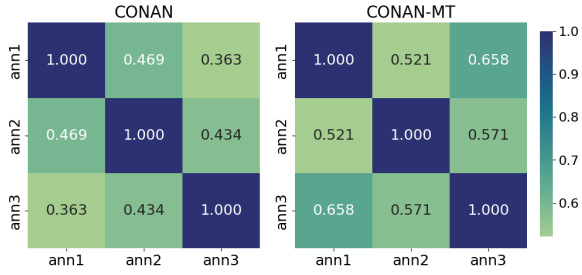


Figure B.3: IAA for the pairwise rank-based evaluation.

Feature	IAA
Relatedness	0.6397
Specificity	0.5947
Richness	0.7806
Coherence	0.8270
Grammaticality	0.7087
Overall	0.6183

Table A1: IAA for feature-wise evaluation.

C.1 Mistral Prompt Format

Provide a brief counter-narrative in response to the user’s hate speech. Ensure the output does not contain line breaks.

###Input:

Hate Speech

###Output:

Counter Narrative

C.2 Mistral-instruct Prompt Format

<s>[INST] Provide a brief counter-narrative in response to the user’s hate speech. Ensure the output does not contain line breaks. *Hate Speech* [/INST] *Counter Narrative*</s>

C.3 Zephyr Prompt Format

<|system|>

Provide a brief counter-narrative in response to the user’s hate speech. Ensure the output does not contain line breaks.</s>

<|user|>

Hate Speech</s>

<|assistant|>

Counter Narrative</s>

C.4 Llama-chat Prompt Format

<s>[INST] «SYS»

Provide a brief counter-narrative in response to the

user’s hate speech. Ensure the output does not contain line breaks.

«/SYS»*Hate Speech* [/INST] *Counter Narrative*</s>

D Correlation Matrix

Correlation between the metrics presented in Section 3.3 and human preference in the CONAN corpus (Table D.4) and in the MT-CONAN corpus (Table D.5).

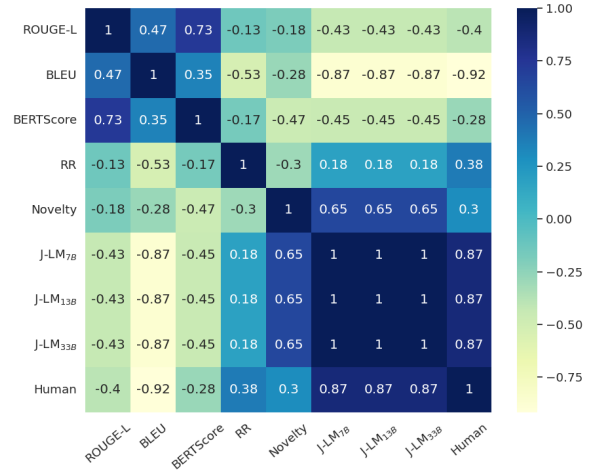


Figure D.4: with the Spearman’s rank correlation between metrics, created using 360 samples from CONAN. The last row of the matrix represents the correlation of all the evaluation methods to human preference.

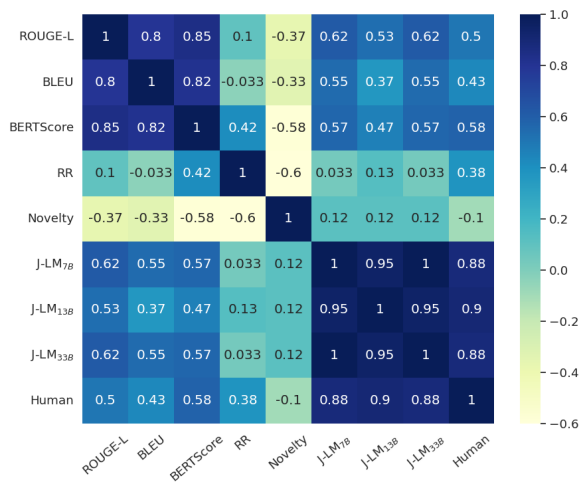


Figure D.5: with the Spearman's rank correlation between metrics, created using 360 samples from MT-CONAN. The last row of the matrix represents the correlation of all the evaluation methods to human preference.