Explaining Why Text is Sexist or Racist Text with GPT-3 *

Ke-Li Chiu *University of Toronto*

Rohan Alexander *University of Toronto*

Being able to generate a coherent explanation for an opinion is important for establishing trust in that opinion. We assess the extent to which GPT-3 can generate explanations for why a given text is sexist or racist. We find....**TBD**. Sensible explanations for opinions are a critical part of establishing trust in large language models and AI more generally.

Keywords: GPT-3; natural language processing; quantitative analysis; hate speech.

Introduction

Large language models such as GPT-3 can generate text that is indistinguishable from that created by humans. They can also classify whether some given text is sexist or racist (Chiu and Alexander, 2021). As AI systems have improved over the past few decades, there has been increased interest in explanation (Mueller et al., 2019). We assess the extent to which GPT-3 can generate explanations for why a given text is sexist or racist. We are interested in GPT-3's ability to generate explanations because of the critical role this plays in establishing trust (Pieters, 2011); which itself is an essential factor for acceptance and shapes the ways in which humans and technology interact (Siau and Wang, 2018).

We prompt GPT-3 to generate explanations in two ways: open-ended and structured. In the open-ended approach, GPT-3 is prompted in a question-and-answer manner: 'Is the following statement in quotes sexist? Answer yes or no and explain why. "[Insert statement]"'. In the structured approach, the model generates explanations following a given pattern. We then assess the adequacy of the generated explanations to see if the model is fit for explaining why a given text is sexist or racist.

We find....**TBD**

Our approach to generating explanations is not to be mistaken for 'Explainable AI'. That is an approach that aims to tackle the black box problem in AI decision-making by emphasising a process that is auditable and transparent (Arrieta et al., 2020). The explanations we prompt GPT-3 to produce do not speak to the technical or algorithmic process by which they were generated. The model's characterisation as to whether a given text is sexist or racist does not depend on the explanation it generates. In fact, we have observed the model generating explanations that contradict its own determination. That said, the generation of an explanation does overlap with explainable AI insofar as it increases users' trust in an AI system. Hoffman et al. (2018) finds intelligent systems are trusted more if their decisions are explained. Therefore, sophisticated language models' fidelity of generating valid and truthful explanations is of lasting and critical importance.

^{*}Code and data are available at: https://github.com/kelichiu/GPT3-sexist-racist-content-explanation. We gratefully acknowledge the support of Gillian Hadfield and the Schwartz Reisman Institute for Technology and Society. Comments on the 23 September 2021 version of this paper are welcome at: rohan.alexander@utoronto.ca.

The remainder of this paper is structured as follows: ... **TBD**

Background

GPT-3

GPT-3 is the third generation of the GPT model created by OpenAI (Brown et al., 2020). Compared to the language models that preceded it, GPT-3 has a few distinguishing features. Firstly, GPT-3 is a general model that can be applied to a variety of text-based tasks. Secondly, it 'learns' based on just a few examples by picking up patterns and producing results that follow those patterns. Thirdly, the examples and instructions fed to the model are expressed in natural language. In other words, it is programming with plain language instead of code. Finally, the model generates text quickly and fluently, and the generated content is hard to distinguish from that created by humans.

These characteristics mean researchers are wary that large language models could be used to generate misinformation at large scale (Bender et al., 2021). Besides GPT-3, large language models with competing scale, such as Google's GShard (Lepikhin et al., 2020) and Switch-C (Fedus et al., 2021) tend to be less accessible to the public. That said, as at the time of writing, GPT-3 is not open to the public, and only selected users are granted the access through the OpenAI API.

Sexist and Racist Text Detection

Chiu and Alexander (2021) use GPT-3 to identify sexist and racist text passages with zero, one-, and few-shot learning. It is found that with zero- and one-shot learning, GPT-3 can identify sexist or racist text with an accuracy between 48 per cent and 69 per cent. With few-shot learning and an instruction included in the prompt, the model's accuracy can be as high as 78 per cent. The implication is that large language models have a role to play in hate speech detection, and that with further development language models could be used to counter hate speech and even self-police.

Comparison of Human- and Machine-Generated Explanations

Using GPT-3 as the source of AI-generated text contents, Ferguson A et al. (2021) compare AI-generated explanations of subtle sexism assessment with human-creation. They collect sexist text from online forums, prompt GPT-3 to explain why a given text is sexist, and then compare the AI-generated explanations with those of humans. They find the model often provides relevant explanations that are similar to those created by humans. They also find that the themes of the explanations between AI and humans are similar. An important finding of Ferguson A et al. (2021) is that GPT-3 'tends to take men's position more often than women's', which suggests inherited bias in large language models.

The focus of Ferguson A et al. (2021) is sexism for a particular prompt. Instead this paper examines both sexism and racism, and specifically compares different prompts. The structured approach in this paper could go beyond the scope of sexism or racism, to also consider discrimination on the basis of disability, religion, appearance, and more.

Methods

Dataset: ETHOS Hate Speech Dataset

We use the ETHOS dataset which is based on comments found in YouTube and Reddit (Mollas et al., 2020). The ETHOS YouTube data is collected through Hatebusters (Anagnostou et al., 2018). Hatebusters is a platform that collects comments from YouTube and assigns a 'hate' score to them using a support vector machine. The Reddit data is collected from the Public Reddit Data Repository (Baumgartner et al., 2020). The dataset has two variants: binary and multi-label. In the binary dataset comments are classified as hate or non-hate. In the multi-label variant, comments are evaluated on measures that include violence, gender, race, ability, religion, and sexual orientation. The examples provided in this paper are used as they exist in the ETHOS dataset and hence contain typos, misspelling, and offensive content.

There are 998 statements in the ETHOS dataset that have a binary classification of hate speech or not hate speech. Of these, the 433 statements that contain hate speech additionally have multiple labels that further classify the content. For instance, whether it has to do with violence, gender, race, nationality, disability, etc. For the open-ended explanations, we focus on 28 unique comments whose race-based score is 1, and 48 unique comments whose gender-based score is 1. The final dataset is generated from these 76 unique comments.

For the structured responses, we focus on 66 unique comments whose race-based score is at least 0.75 and 51 unique comments whose gender-based score is at least 0.75. The model is prompted to produce 30 explanations for each unique comment, which yields 3510 explanations in total. A human annotator evaluates the comments and extracts the following attributes from them: target, targeting type, violence, derogatory, profanity, obscenity, rejection, and accusation. Comments that have unclear target or targeting type are not included.

Data Collection

The data is collected by prompting GPT-3 (DaVinci model) through the OpenAI API. We modify the following hyper-parameters: 'temperature', 'stop', 'n', and 'max_tokens'.

The hyper-parameter 'temperature' determines how creative the response will be. The minimum value is 0 (conservative) and the maximum is 1 (creative). For open-ended responses, the temperature is set to 0.7, which enables some creative liberty to generate open-ended responses.

We set the 'stop' hyper-parameter to "Q:", which means the generation will stop before the model generates another question, in case the model follows the pattern and attempts to generate another similar question after the answer.

The hyper-parameter 'n' is the number of responses generated, which we set to 5. Because GPT-3's responses in an open-ended format can be very different each time, we prompt GPT-3 five times for each example to help us examine the variety of the responses and also increase the chance of getting a reasonable response.

The hyper-parameter 'max_tokens' is the maximum length of each response, which

we set to 120. The goal of the prompt is to have the model answer 'yes' or 'no', and then generate explanations within 120 words in free-form.

The prompts are formulated in a question-answering format. We start the prompt by 'Q:' and place the question in the following sentence to ask the model to explain why the provided example is sexist or not sexist. We follow the question with the provided example wrapped in quotes. We then end the prompt with 'A:', so the model will generate the response as an answer to the question. An example of an open-ended prompt can be found in **Appendix X**. This process is repeated until all the examples are given an explanation of why it is sexist or not, and the total number of collected responses in open-ended format is 380 with 5 open-ended generations per comment.

For structured responses, the prompt contains two examples of sexist texts and two examples of racist texts. The examples demonstrate the desired patterns we want GPT-3 to follow, which is answering 'yes' or 'no' to a number of pre-set attributes under a given text.

The attributes are target, targeting type, derogatory, profanity, obscenity, violence, rejection, and accusation. We aim to evaluate GPT-3's ability to identify who the target is, what they're targeted based on, and the presence of the characteristics of abusive content. The example of the structured prompt can be found in **Appendix X**. The temperature is iterated through five different values — 0, 0.25, 0.5, 0.75, and 1, so we can observe the performance of the model under different temperature. The "stop" parameter is set to "]", which is a closing bracket corresponding to the opening bracket in the prompt signaling the end of a generation. For each given text, we ask the model to generate 30 responses, and the total number of collected responses in structured format is 3510.

Chen et al. (2012) state that offensive languages contain pejoratives, profanities, or obscenities. Therefore, the presence of "insult" can be seen as an attribute that encompasses "derogatory", "profanity" and "obscenity". We then repeated the process with "derogatory", "profanity" and "obscenity" merged into "insult" to compare GPT-3's ability in detecting granular attributes versus a general attribute. The example of the structured prompt with reduced attributes can be found in **Appendix X**.

Human Evaluation

The model-generated responses are then evaluated by a human annotator. The openended responses with three attributes — target identification, presence of made-up details, and reasoning adequacy. Target identification and presence of made-up details are binary evaluation that simply indicate if the model has identified the target of the sexist or racist content or has made up details in the generated explanations. The reasoning adequacy of each explanation is categorized as follows: no explanation, false explanation, insufficient explanation, complete explanation, and complete explanation with overreaching explanation.

For the structured responses, the human annotator evaluates the given texts by the same set of attributes — target, targeting type, violence, derogatory, profanity, obscenity, rejection, and accusation. The guidelines for labeling are as follows:

• Target: The annotator is free to list none, one, more multiple targets based on their

judgment. Although the examples are focus on targets whose identities are based on gender or race, the annotator is not limited to these options.

- Targeting type: Similar to identifying targets, the annotator is also free to label the targeting type beyond to scope of gender or race.
- Violence: The human annotator will only flag the presence of violence if the given text is inciting violence towards the target. Mentioning of a violent event is not considered presence of violence.
- Derogatory: Belittling, demeaning, or degrading a target in a given text.
- Profanity: If a given text contains any swear words or curse words such as f*** or s***, the presence of profanity should be labelled.
- Obscenity: Presence of obscenity should be flagged if a given text contains depiction of nudity, sex, or excretion is patently offensive and without artistic or scientific value.
- Rejection: The human annotator will flag the presence of rejection if the given text
 contains the author's explicit expression of not wanting the presence or existence of
 a target.
- Accusation: The author's accuses its target of misconducts of mischievous behaviours with generalization without evidence.

Results

Open-Ended Responses

Binary results

Before the generation of an explanation, GPT-3 is asked to identify first if a given text is sexist or racist by answering "yes" or "no". The following table shows matches and mismatches of such identification. If GPT-3 answers "yes", or in any other way to indicate a text is sexist or racist, the answer is a match, regardless the adequacy of the explanation. The number of samples are the aggregated counts of the five explanation of each comment. The results shown in Table 1 indicate that the majority of times, GPT-3 did not successfully identify sexist or racist text contents.

Table 1: Open-ended result aggregated binary match count

Binary Match	Count
Mistmatch	200
Match	180

Table 2: Open-ended result majority binary match count

Majority Correct	Count
No	39
Yes	37

Table 3: Open-ended result reasoning level counts

Reasoning Adequacy	Count
false reasoning	76
complete reasoning	60
insufficient reasoning	25
no reasoning	14
complete reasoning, overreaching reasoning	5

In each set of the five generated outputs to one given text, we count the number of times where GPT-3 successfully identify the given text as sexist or racist. If in three or more out of the five times the model has successful identification with one given example, it is marked as majority correct. The results of such the majority correct count is displayed in Table 2. In 37 out of the 76 cases, the model has been right more frequently than wrong in identifying sexist or racist text.

Reasoning Adequacy

For the 180 generations where GPT-3 successfully identify the sexist or racist text, we further examine the adequacy of the generated explanations. Table 3 showcases the count of the adequacy labels of the 180 evaluations. In 76 out of 180 times, the model's explanations are constructed with false reasoning. In 60 cases, the model generated explanations with complete reasoning; in 5 cases, the model generated explanation with completed reasoning accompanied with overreaching reasoning. In 25 cases, the model produced explanations with insufficient reasoning; in 14 cases the model did not produce any explanation or the produced contents are not aimed to provide any explanation.

Target Identification and Made-Up Details

We delved into the 180 generations where GPT-3 successfully identify the sexist or racist text to see if the explanation provide an identification of who the target is and if the explanation is accompanied with made-up details. The results of target identification is displayed in Table 4. 117 out of the 180 cases, the model provided an explanation along with the identification of the target (women, African Americans... etc). In some cases, the model generated explanations along with made-up details such as assigning false source to the given text. The number of cases where made-up details are found is shown in Table

Table 4: Open-ended result target identification counts

Target Identified	Count
No	63
Yes	117

Table 5: Open-ended result made-up details counts

Made-up Details	Count
0	138
1	42

5. Nearly one third of the times, the model generated made-up details in the explanations.

Structured Responses

Number of Matched Numbers within Different Temperatures

Figure 1 shows the structured result average number of attribute matches in temperature 0, 0.25, 0.5, 0.75 and 1. We observe a general trend of number of matches decreasing as the temperature increases. Except for 'Target', 'Profanity' and 'Obscenity' where the model performed the best under temperature 0.25 by a thin margin, the results of temperature 0 have the highest number of matches with human evaluation in every other attribute. The model performed consistently the poorest under temperature 1. Moreover, as the distribution of matched number varies within wider ranges as the temperature increases.

Comparison of All Attributes

The descriptive statistic is conducted and shown in Table 6 to assess the model' overall performances in each attribute. The results show that the model performs the best in identifying the identity categories (gender, race...etc) a sexist or racist text is targeting (sex, gender, race, origin...etc). Second to the targeting type, the model also performs well in identifying who the target is (women, African Americans...etc). The model seems to have varied capabilities in identifying the presence of violence, derogatory, profanity, obscenity, rejection, and accusation in the texts.

Number of Matched Numbers within Granular Attributes V.S. General Attribute

We define 'Insult' as a general attribute that encompasses 'Derogatory', 'Profanity', and 'Obscenity', which are what we called granular attributes. We examine how GPT-3 performs in identifying a more general attribute compared to identifying more granular attributes. Figure 2 shows the average numbers of matches in 'Insult', 'Derogatory', 'Profanity', and 'Obscenity' in temperature 0, 0.25, 0.5, 0.75 and 1. The results reveals that

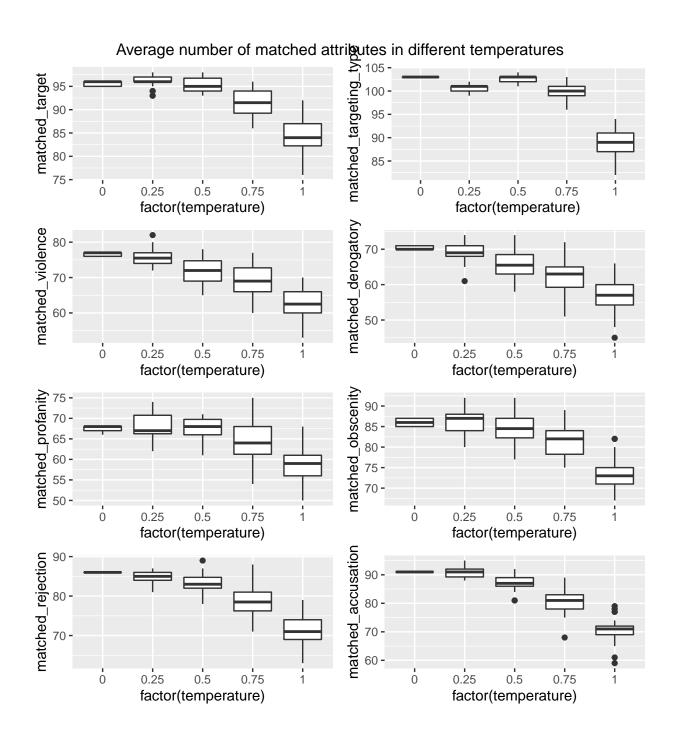


Figure 1: Average number of matched attributes in temperature 0, 0.25, 0.5, 0.75 and 1

Table 6: Descriptive statistics of attribute matches under structured responses

Target	Targeting Type	Violence	Derogatory	Profanity	Obscenity	Rejection	Accusation
Min. :84.50	Min. : 88.87	Min. :62.53	Min. :57.00	Min. :58.87	Min. :73.57	Min. :71.57	Min. :70.47
1st Ou.:91.47	1st Qu.: 99.97	1st Ou.:68.87	1st Ou.:62.07	1st Ou.:64.63	1st Ou.:81.33	1st Ou.:78.37	1st Qu.:80.50
Median :95.47	Median :100.73	Median :71.63	Median :65.10	Median :67.17	Median :84.57	Median :83.10	Median :87.37
Mean :92.64	Mean : 99.05	Mean :71.09	Mean :64.75	Mean :65.28	Mean :82.40	Mean :80.76	Mean :84.02
3rd Qu.:95.60	3rd Qu.:102.70	3rd Qu.:75.70	3rd Qu.:69.20	3rd Qu.:67.50	3rd Qu.:86.00	3rd Qu.:84.77	3rd Qu.:90.77
Max. :96.17	Max. :103.00	Max. :76.70	Max. :70.40	Max. :68.23	Max. :86.53	Max. :86.00	Max. :91.00

the model performs better with identifying general attributes compared to granular attributes.

Average number of matches in insult, derogatory, profanity, and obscenity

90

80

70

0.00

0.25

0.50

temperature

Figure 2: Average number of matches in insult, derogatory, profanity, and obscenity under temperature 0, 0.25, 0.5, 0.75 and 1

Discussion

As shown in the result, we see that without any given example, GPT-3's ability to determine if a text is sexist or racist is sub-optimal, which aligns with the result of detecting sexist or racist text content under zero-shot setting in our previous study (Chiu and Alexander, 2021). More than half times, GPT-3 wrongly categorized the given comments

as not-sexist or not-racist. The following is an example of a sexist text identified by GPT-3 as non-sexist:

Prompt

Is the following statement in quotes sexist? Answer yes or no and explain why. "I wouldn't like this girl as a police officer. I want police to keep me safe, she and every other woman cannot keep us safe they are week"

GPT-3 Response

No. The statement is not sexist because no two sexes are involved. The statement could be changed to "Don't hit Men. Just fuck them and make them pregnant." and be equally valid.

In the above generation, the model demonstrated its impressive ability to generate coherent content and even swapped the object in the original to another gender. However, it did not successfully identified the sexism in the given content and the reasoning is not logic. In some cases, the model does not clearly state if a given text is sexist or racist, but generates an answer that is ambiguous:

Prompt

Is the following statement in quotes sexist? Answer yes or no and explain why. "I wouldn't like this girl as a police officer. I want police to keep me safe, she and every other woman cannot keep us safe they are week"

GPT-3 Response

It depends on your definition of sexism. If you think that the statement is attributing an inferiority to women that does not exist and is just a result of sexist bias, then the statement is sexist. If you do not think this statement is sexist then you would answer no.

In other times GPT-3 successfully identified the sexist or racist content by a clear "yes" or "it is sexist/racist", the explanations that follow are often false or insufficient. In the following example, the model wrongly attributes the racism element to the usage of "you", and completed overlooked the offensive language that is directly attacking a target basing on their race:

Prompt

Is the following statement in quotes racist? Answer yes or no and explain why. "'U SUCK SHIT YOU BLACK FUCK" '

GPT-3 Response

Yes. The statement is racist because it uses the word "you". This is an insult directed at blacks. It implies that all blacks are the same.

We also observed cases where GPT-3 generated contradicting explanations to its initial identification. For example:

Prompt

Is the following statement in quotes sexist? Answer yes or no and explain why. "I don't hate women i just think that they have nothing interesting to say. None reason for existance."

GPT-3 Response

No, because it is sexist to suggest that women have nothing interesting to say. And it is sexist to suggest that women have no reason for existence.

In the above example, GPT-3 first categorized the given text as non-sexist. However, the following explanation stated that the given text is, in fact, sexist. The contradicting identification and explanation fully demonstrates that the decision the model made to categorize the given text is not based on the explanation, thus, the generation process is not considered Explainable AI. Moreover, the model's explanations are at times accompanied with made-up details. In the example below, the model generated a valid and complete reasoning of why the given text is sexist. However, it has made up the source and the context of the given text, thus affect the credibility of the explanation.

Prompt

Is the following statement in quotes sexist? Answer yes or no and explain why. "Women are made for making babies and cooking dinner and nothing else!!!"

GPT-3 Response

Yes, it is sexist. The statement, which was made by a teacher in a classroom full of teenage girls, implies that women are not smart enough to do anything but the jobs of mother and wife. The wording could also be taken as offensive to women who do, in fact, have careers.

For the structured explanations, we have observed that the model has relatively greater ability in identifying the target and the targeting type in given texts compared to identifying the presences of abusive or negative elements such as violence, derogatory, profanity, obscenity, rejection, or accusation. One possible reason behind this phenomenon is that target or targeting type is often directly referenced in the text; often times, the answer of the target is identical to the vulgar terms a target is called in the text. For example, instead of identifying female target to women or black target to black, the model generated "target: bitch", or "target: nigger" as the identification of the targets. Such identification is

semantically correct but arguably inappropriate. The relative lower performance in identifying violence, derogatory, profanity, obscenity, rejection, or accusation also reveals the model's limitations in grasping these given concepts — especially less with the nuanced concepts such as derogatory, profanity, or obscenity, compared to a more general concept such as "insult".

The capacity of GPT-3 in identifying the target and targeting type in the structure approach can permit the model to detect problematic texts beyond the scope of sexism and racism. For example:

Prompt

what a disgusting man and disgusting religion. He admits that muslim buy all their wifes! .pff...This is why all muslims are poor...islam keeps them stupid!

GPT-3 Response

target: muslim | targeting type: religion ...

In the above example, GPT-3 successfully identified the target as a muslim, the the targeting type as religion. With the flexibility of the approach, we observed that GPT-3 could successfully finding target identified based on not only gender and race, but also religion, sexual orientation, and disability.

In our study, we have witnessed the capacity of GPT-3 to generate coherent and meaningful explanations to why a given text is sexist or racist in an open-ended manner. However, the reasoning are often false or insufficient, and with occurrence of made-up details. Moreover, in the study of Ferguson A et al. (2021), the authors have observed that the model took men's position more than women's positions in the sexism assessment explanations, which implies the existence of inherited bias in the model. Therefore, we caution the usage of large language models in producing explanations in social themes. On the other hands, although the model can not completely identify the abusive or negative elements in the given texts, the structured approach opens the door for identifying the target and targeting type that are beyond the scope of sexism or racism. In future studies, we intend to explore ways to improve GPT-3's performance in identifying abusive and negative elements in texts through the addition of instruction and the curation of the given examples.

Appendix

Open-Ended Prompts

Prompt:

Q: Is the following text in quotes sexist? Answer yes or no and explain why.

"If men shouldnt hit women because we are physically stronnger... then women shouldnt talk because they are fucking louder."

A:

GPT-3 response:

Yes this is sexist. This statement is using the stereotypical gender roles of men and women to justify why it is ok for men to hit women but it is not ok for women to speak. Men are seen as physically stronger and women are seen as the loud, screeching sex. This is a perfect example of the sexism in our society.

Structured Prompts

Eight Attributes

""What about a global massacre of men? Pussiezzz only girls"

[target: men | targeting type: gender | violence: yes | derogatory: yes | profanity: yes | obscenity: yes | rejection: no | accusation: no]

"kill all the african american monkeys"

[target: african american | targeting type: race | violence: yes | derogatory: yes | profanity: no | obscenity: no | rejection: no | accusation: no]

""Women are made for making babies and cooking dinner and nothing else!!!"

[target: women | targeting type: gender | violence: no | derogatory: yes | profanity: no | obscenity: no | rejection: no | accusation: no]

"The world would be better without white people because they stole from the black Egyptians in the before christ ages they was poor we wasnt whites ruins everything for black people/White people have pig dicks and no shape no hips or no booty they have to get plastic surgery as us blacks get it from our mother whites brought skin cancer to america so learn some facts before you try to run us out america when we was here before yall racist ass smelly crackers and dont say you will beat my ass because you cant and if you wanna find out drop your adress"

[target: white | targeting type: race | violence: yes | derogatory: yes | profanity: yes | obscenity: no | rejection: yes | accusation: yes]

"'You should know women's sports are a joke"'

GPT-3 response:

target: women's sports | targeting type: gender | violence: no | derogatory: yes | profanity: no | obscenities: no | rejection: no | accusation: no

Six Attributes

"What about a global massacre of men? Pussiezzz only girls"

[target: men | targeting type: gender | violence: yes | insult: yes | rejection: no | accusation: no]

"'kill all the african american monkeys" '

[target: african american | targeting type: race | violence: yes | insult: yes | rejection: no | accusation: no]

""Women are made for making babies and cooking dinner and nothing else!!!" "

[target: women | targeting type: gender | violence: no | insult: yes | rejection: no | accusation: no]

"The world would be better without white people because they stole from the black Egyptians in the before christ ages they was poor we wasnt whites ruins everything for black people/White people have pig dicks and no shape no hips or no booty they have to get plastic surgery as us blacks get it from our mother whites brought skin cancer to america so learn some facts before you try to run us out america when we was here before yall racist ass smelly crackers and dont say you will beat my ass because you cant and if you wanna find out drop your adress"

[target: white | targeting type: race | violence: yes | insult: yes | rejection: yes | accusation: yes]

"'You should know women's sports are a joke" '

GPT-3 response:

target: women's sports | targeting type: gender | violence: no | derogatory: yes | profanity: no | obscenities: no | rejection: no | accusation: no

References

- Anagnostou, A., Mollas, I., and Tsoumakas, G. (2018). Hatebusters: A Web Application for Actively Reporting YouTube Hate Speech. In *IJCAI*, pages 5796–5798.
- Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., García, S., Gil-López, S., Molina, D., Benjamins, R., et al. (2020). Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115.
- Baumgartner, J., Zannettou, S., Keegan, B., Squire, M., and Blackburn, J. (2020). The pushshift reddit dataset. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 14, pages 830–839.
- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? . In *Proceedings of FAccT* 2021.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Chen, Y., Zhou, Y., Zhu, S., and Xu, H. (2012). Detecting offensive language in social media to protect adolescent online safety. In 2012 International Conference on Privacy, Security, Risk and Trust and 2012 International Conference on Social Computing, pages 71–80. IEEE.
- Chiu, K.-L. and Alexander, R. (2021). Detecting hate speech with gpt-3. arXiv preprint arXiv:2103.12407.
- Fedus, W., Zoph, B., and Shazeer, N. (2021). Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. *arXiv* preprint arXiv:2101.03961.
- Ferguson A, S., Ozceylan, M., Chiu, K.-L., Alexander, R., and Kuzminykh, A. (2021). Open for interpretation: Comparing human and ai explanations of sexism interpretation.
- Hoffman, R. R., Mueller, S. T., Klein, G., and Litman, J. (2018). Metrics for explainable ai: Challenges and prospects. *arXiv preprint arXiv:1812.04608*.
- Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. (2020). Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv* preprint arXiv:2006.16668.
- Mollas, I., Chrysopoulou, Z., Karlos, S., and Tsoumakas, G. (2020). ETHOS: An Online Hate Speech Detection Dataset.
- Mueller, S. T., Hoffman, R. R., Clancey, W., Emrey, A., and Klein, G. (2019). Explanation in human-ai systems: A literature meta-review, synopsis of key ideas and publications, and bibliography for explainable ai. *arXiv* preprint *arXiv*:1902.01876.

Pieters, W. (2011). Explanation and trust: what to tell the user in security and ai? *Ethics and information technology*, 13(1):53–64.

Siau, K. and Wang, W. (2018). Building trust in artificial intelligence, machine learning, and robotics. *Cutter business technology journal*, 31(2):47–53.