

Detecting Hate Speech with GPT-3 *

Ke-Li Chiu *University of Toronto*

Rohan Alexander *University of Toronto and Schwartz Reisman Institute*

Sophisticated language models such as OpenAI’s GPT-3 can generate hateful text that targets marginalized groups. Given this capacity, we are interested in whether large language models can be used to identify hate speech and classify text as sexist or racist? We use GPT-3 to identify sexist and racist text passages with zero-, one-, and few-shot learning. We find 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. We conclude 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.

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

Introduction

This paper contains language and themes that are offensive.

Natural language processing (NLP) models focus on using normal words, often written text, as their data. For instance, one might have content from a large number of books and want to group them into themes. Sophisticated NLP models are being increasingly embedded in society. For instance, Google Search uses an NLP model called BERT to better understand what is meant by a word given its context. Some sophisticated natural language processing (NLP) models, such as OpenAI’s GPT-3, can additionally produce text as an output.

The text produced by sophisticated NLP models can be hateful. In particular, there have been many examples of text being generated that targets marginalized groups based on their sex, race, sexual orientation, and other characteristics. For instance, ‘Tay’ was a chatbot that was released by Microsoft in 2016 on Twitter. Within hours of being released, some of its tweets were sexist. Large language models are trained on enormous datasets from various sources. This means that untruthful statements, human biases, and abusive language are inevitably included. Further, even when the models do not possess intent, they risk producing synthetic texts that are offensive or discriminatory and thus cause unpleasant, or even triggering, interaction experiences (Bender et al., 2021).

*Code and data are available at: <https://github.com/kelichiu/GPT3-hate-speech-detection>. We gratefully acknowledge the support of Gillian Hadfield, the Schwartz Reisman Institute for Technology and Society, and OpenAI for providing access to GPT-3 under the academic access program. We thank Amy Farrow, Haoluan Chen, John Giorgi, Mauricio Vargas Sepúlveda, Monica Alexander, Noam Kolt, and Tom Davidson for helpful discussions and suggestions. Please note that we have added asterisks to racial slurs and other offensive content in this paper, however the inputs and outputs did not have these. Comments on the 21 February 2022 version of this paper are welcome at: rohan.alexander@utoronto.ca.

Often the datasets that underpin these models consist of, essentially, the whole public internet. This source raises concerns around three issues: exclusion, over-generalization, and exposure (Hovy and Spruit, 2016). Exclusion happens due to the demographic bias in the dataset. In the case of language models that are trained on US and UK English scraped from the Internet, datasets may be disproportionately white, male, and young. Therefore, it is not surprising to see white supremacist, misogynistic, and ageist content being over-represented in training datasets (Bender et al., 2021). Over-generalization stems from the assumption that what we see in the dataset represents what occurs in the real world. Words such as ‘always’, ‘never’, ‘everybody’, or ‘nobody’ are frequently used for rhetorical purpose instead of their literal meanings. However, language models do not always recognize this and make inference based on generalized statements using these words. For instance, hate speech commonly uses generalized language for targeting a group such as ‘all’ and ‘every’, and a model trained on these statements may generate similarly generalized and harmful statements. Finally, exposure refers to the relative attention, and hence considerations of importance, given to something. In the context of NLP this may be reflected in the emphasis on English-language created under particular circumstances, rather than another languages or circumstances that may be more prevalent.

While these, and other, issues give us pause, the dual-use problem, which is that the same technology can be applied to both good and bad uses, provides motivation. For instance, while stylometric analysis can reveal the identity of political dissenters, it can also solve the unknown authorship of historic text (Hovy and Spruit, 2016). In this paper we are interested in whether given large language models can produce harmful language, can they also identify, or learn to identify, harmful language?

Even though large NLP models do not have a real understanding of language, the vocabularies and the construction patterns of harmful languages can be thought of as known to them. We show that this knowledge can be used to identify abusive language and even hate speech. In particular we consider 120 different extracts that have been categorized as ‘racist’, ‘sexist’, or ‘neither’ in the single-category settings and 240 different extracts in the mixed-category settings. We ask GPT-3 to classify these based on zero-, one-, and few-shot learning, with and without instruction. We find that the model performs best with few-shot learning when an instruction is included. In that setting the model can accurately classify around 78 per cent of the extracts. If language models can be used to identify abusive language, then not only is there potential for them to counter the production of abusive language by humans, but they could also potentially self-police.

Background

Language models, Transformers and GPT-3

In its simplest form, a language model involves assigning a probability to a certain sequence of words. For instance, the sequence ‘the cat in the hat’ is probably more likely than ‘the cat in the computer’. We typically talk of tokens, or collections of characters, rather than words, and a sequence of tokens constitutes different linguistic units — words, sentences, and even documents (Bengio et al., 2003). Language models predict

the next token based on inputs. If we consider each token in a vocabulary as a dimension, then the dimensionality of language quickly becomes large (Rosenfeld, 2000). Over time a variety of statistical language models have been created to nonetheless enable prediction. The n-gram is one of the earliest language models. It works by considering the co-occurrence of tokens in a sequence. For instance, given the four-word sequence, ‘the cat in the’, it is more likely that the fifth word is ‘hat’ rather than ‘computer’. In the early 2000s, language models based on neural networks were developed, for instance Bengio et al. (2003). These were then built on by word embeddings language models in the 2010s in which the distance between tokens represents how related those tokens are, for instance Turian et al. (2010). In 2017, Vaswani et al. (2017) introduced the Transformer, which marked a new era for language models. The Transformer is a network architecture for neural networks that can be trained more quickly than many other approaches (Vaswani et al., 2017). Now most representative pre-trained language models, such as Google’s model: Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2018), as well as OpenAI’s: Generative Pre-trained Transformer (GPT)-2 (Radford et al., 2019), and GPT-3 (Brown et al., 2020), are built on this architecture. These models are widely used, for instance BERT is used by Google search.

GPT-3 is the third generation of the Generative Pre-trained Transformer models created by OpenAI. OpenAI is a private company in California that develops artificial intelligence models. Until January 11, when Fedus et al. (2021) announced a Transformer model with a trillion parameters, GPT-3 was the largest, publicly-known, Transformer language model. GPT-3 is distinctive from its predecessors because of few-shot learning. This means that GPT-3 can ‘learn’ to perform a new task based on only a few examples, expressed in natural language, instead of a fine-tuning process that can require a large amount of data. GPT-3 has led to unexpected NLP applications, such as computational code generation given natural language prompts.

Like other language models, GPT-3 has also generated inappropriate or even hateful content. For instance, McGuffie and Newhouse (2020) demonstrated the use of GPT-3 in mass-producing radicalized text targeting the Islamic population. And Lin et al. (2021) show that GPT-3 and similar language models can propagate misconceptions that could deceive human readers. For instance, when asked ‘Who really caused 9/11?’, they found that GPT-3 provided the false statement ‘The US government caused 9/11.’

Hate speech detection

There is no commonly held definition of hate speech and different legal jurisdictions have different definitions, as do different companies and other groups. One definition is ‘the intentional verbalization of prejudice against a social group’ (Kennedy et al., 2018). Detecting hate speech is difficult because the definition of hate speech varies depending on the complex intersection of the topic of the assertion, the context, the timing of the post, synchronized world events, and the identity of speaker and recipient (Schmidt and Wiegand, 2017). Moreover, it is difficult to discern hate speech from merely offensive language (Davidson et al., 2017). Hate speech detection is of interest to academic researchers in a variety of domains including computer science (Srba et al., 2021) and sociology (Davidson et al., 2017). It is also of interest to industry, for instance to maintain standards on social

networks, and in the judiciary to help identify and prosecute crimes. Since hate speech is prohibited in several countries, misclassification of hate speech can become a legal problem. For instance, in Canada, speech that contains ‘public incitement of hatred’ or ‘wilful promotion of hatred’ is specified by the Criminal Code ([Criminal Code, 1985](#)). Policies toward hate speech are more detailed in some social media platforms. For instance, the Twitter Hateful Conduct Policy states:

You may not promote violence against or directly attack or threaten other people on the basis of race, ethnicity, national origin, caste, sexual orientation, gender, gender identity, religious affiliation, age, disability, or serious disease. We also do not allow accounts whose primary purpose is inciting harm towards others on the basis of these categories.

[Twitter \(2021\)](#)

There has been a large amount of research focused on detecting hate speech. And as part of this various hate speech datasets have been created and examined. For instance, [Waseem and Hovy \(2016\)](#) create a dataset that captures hate speech in the form of racist and sexist language that includes domain expert annotation. They use Twitter data, and annotate 16,914 tweets: 3,383 as sexist, 1,972 as racist, and 11,559 as neither. They had a high degree of annotator agreement, and most of the disagreements were to do with sexism, and often explained by an annotator lacking apparent context. And, [Davidson et al. \(2017\)](#) trains a classifier to distinguish between hate speech and offensive language. To define hate speech, they use online ‘hate speech lexicon containing words and phrases identified by internet users as hate speech’. It is important to note that even these datasets have bias. For instance, [Davidson et al. \(2019\)](#) found racial bias in five different sets of Twitter data annotated for hate speech and abusive language. They found that tweets written in African American English are more likely to be labeled as abusive.

Methods

We examine the ability of GPT-3 to identify hate speech in zero-shot, one-shot, and few-shot settings. There are a variety of parameters, such as temperature, that control the degree of text variation. Temperature is a hyper-parameter between zero and one. Lower temperatures mean that the model places more weight on higher-probability tokens. To enhance consistency, the temperature is set to zero in our experiments and not varied. There are two categories of hate speech that are of interest in this paper. The first targets the race of the recipient, and the other targets the gender of the recipient. With zero- and one-shot learning, the model identifies hate speech one category at a time. With few-shot learning, the categories are mixed, and the model is asked to classify an input as sexist, racist, or neither. Zero-shot learning means an example is not provided in the prompt. One-shot learning means that one example is provided, and few-shot means that two or more examples are provided.

Dataset

We use the ETHOS dataset created by [Mollas et al. \(2020\)](#). ETHOS is based on comments found in YouTube and Reddit. 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. That hate-score is only used to decide whether to consider the comment further or not. The Reddit data is collected from the Public Reddit Data Repository ([Baumgartner et al., 2020](#)). The classification is done by contributors to a crowd-sourcing platform. They are first asked with an example contains hate speech, and then, if it does, whether it incites violence and other additional details. The dataset has two variants: binary and multi-label. In the binary dataset comments are classified as hate or non-hate based. In the multi-label variant, the comments are evaluated on measures that include violence, gender, race, ability, religion, and sexual orientation. The examples provided in this paper are from the ETHOS dataset and hence contain typos, misspelling, and offensive content.

We use all of the 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 classify the content, for instance, does the comment have to do with violence, gender, race, nationality, disability, etc. We considered all of the 136 statements that contain race-based hate speech, and we focus on the 76 whose race-based score is at least 0.5. Similarly, we considered the all of 174 statements that contain gender-based hate speech, and again focused on the 84 whose gender-based score is at least 0.5. To create a balanced dataset we randomly draw 120 of these statements, using stratified sampling, to result in 60 sexist statements and 60 racist statements. For comparison purposes, we randomly draw 120 statements that are classified as not hate speech. The balanced sample enables us to more easily understand our results.

Zero-shot learning

With zero-shot learning, the model is not given an example and is expected to generate an answer without additional information. In this setting, GPT-3 is asked to identify one passage of text, in one category, each time. A simple instruction is given asking the model to answer only ‘Yes’ or ‘No’. We ask the model to perform the identification of sexist and racist comments 120 times. The model considers 30 sexist comments, 30 racist comments, and 60 comments that are neither sexist nor racist. An example of zero-shot learning for recognizing a comment that is labeled as sexist is:

Prompt:

Is the following text sexist? Answer yes or no.

‘The thing is women are not equal to us men and their place is the home and kitchen’

Answer:

GPT-3 response:

Yes.

The model provides this answer based on the dataset that it was trained on.

One-shot learning

With one-shot learning, similarly to the zero-shot learning case, the model is asked to identify one passage of text in one category at a time. In the one-shot learning case, the model is given one example from a category and is asked to identify whether the next passage of text belongs to the selected category. Again, a simple instruction is also given to ask the model to answer only 'Yes' or 'No', and the same break-down of 30 sexist comments, 30 racist comments, and 60 comments that are neither sexist nor racist is considered. That is, the model performs this task 120 times with different comments each time. An example of one-shot learning for recognizing a comment that is labeled as sexist is:

Prompt:

The following text in quotes is sexist:

'Feminism is a very terrible disease'

Is the following text sexist? Answer yes or no.

'She is heavily relying on him to turn the other cheek... tough talking demon infested woman.'

Answer:

GPT-3 response:

Yes.

Few-shot learning

Under the few-shot setting, three examples are provided to the model. There are two settings considered in the few-shot experiments: single-category, and mixed-category. In the single-category setting, the model is given two examples that belong to the selected category, and one example that is not hate speech. The answer is expected to be binary. For instance, the model is expected to categorize a passage of text as 'sexist' or 'not-sexist'. An example of this is:

Prompt:

'Too bad women don't know how to kill themselves': sexist.

'You should use your time to arrest murderers not little kids': not-sexist.

'Now they know better than this s*** lol they dudes. The stronger sex. The man supremacy': sexist.

'The thing is women are not equal to us men and their place is the home and kitchen:'

GPT-3 response:

sexist.

The assignment of sexist or not in the prompt is based on the ETHOS dataset classification.

On the other hand, in the mixed model, the model is given an example of sexist text, an example of racist text, and an example of text that is not hate speech. In the mixed categories setting, the model is expected to categorize the given passage of text as 'sexist', 'racist', or 'neither'. There are ten sets of examples; each set of examples is used for the model to classify 60 racist comments, 60 sexist comments, and 120 comments that are in neither category. This results in 2,400 responses in total. The ten sets of examples are available in Table 16 in the Appendix. An example of this is:

Prompt:

'That girl still need a good ass whooping, stupid ass b****!': sexist.

'U SUCK S*** YOU BLACK F***': racist.

'You deserve way more subscribers. The s*** you say is f***ing gold.': neither.

'Every girl on this planet curses her mother for their birth. They know that they are the powerless sex and they pray that they had been born as men! With d***s and stuff.':

GPT-3 response:

sexist.

Finally, in the mixed model, we added a brief instruction to restrict the generated answers to 'racist', 'sexist' or 'neither'. This is to prevent the model from generating answers that are out of scope. An example of this is:

Prompt:

Classify the following texts into 'racist', 'sexist' or 'neither'.

'That girl still need a good ass whooping, stupid ass b****!': sexist.

'U SUCK S*** YOU BLACK F***': racist.

'You deserve way more subscribers. The s*** you say is f***ing gold.': neither.

'Every girl on this planet curses her mother for their birth. They know that they are the powerless sex and they pray that they had been born as men! With d***s and stuff.':

GPT-3 response:

sexist.

Table 1: Classification of racist statements with zero-shot learning

Actual classification	GPT-3 classification	
	Not racist	Racist
Not racist	0	30
Racist	7	23

Table 2: Classification of sexist statements with zero-shot learning

Actual classification	GPT-3 classification	
	Not sexist	Sexist
Not sexist	5	25
Sexist	0	30

Result

Zero-shot learning

The results of the zero-shot experiments are presented in Table ?? . The model has 35 matches and 25 mismatches in the sexist category, and 23 matches and 37 mismatches in the racist category. The model performs better when identifying sexist comments compared with identifying racist comments. However, the overall ratio of matches and mismatches is 58:62. In other words, the accuracy in identifying hate speech in the zero-shot setting is 48 per cent.

One-shot learning

The results of the one-shot learning experiments are presented in Table ?? . The model has 46 matches and 14 mismatches in the racist category, and 37 matches and 23 mismatches in the sexist category. In contrast with the result generated from zero-shot learning, the model performs slightly better in identifying racist comments compared with identifying sexist comments. The general performance in the one-shot setting is also better than in the

Table 3: Classification of hate speech with zero-shot learning

Actual classification	GPT-3 classification	
	Not hate speech	Hate speech
Not hate speech	5	55
Hate speech	7	53

Table 4: Performance of zero-shot learning in hate speech classification

	Percent
Racism	
Accuracy	38.33
Precision	43.40
Recall	76.67
F1	55.43
Sexism	
Accuracy	58.33
Precision	54.55
Recall	100.00
F1	70.59
Overall	
Accuracy	48.33
Precision	49.07
Recall	88.33
F1	71.43

Table 5: Classification of racist statements with one-shot learning

Actual classification	GPT-3 classification	
	Not racist	Racist
Not racist	21	9
Racist	5	25

zero-shot setting. The overall ratio of matches and mismatches is 83:37. In other words, the accuracy of identifying hate speech in the one-shot setting is 69.2 per cent.

Few-shot learning – single category

The results of the single-category, few-shot learning, experiments are presented in Table ???. The model has 41 matches and 19 mismatches in the racist category, and 42 matches and 18 mismatches in the sexist category. Here, the model performs almost equally well at identifying racist comments compared with identifying sexist comments. The general performance in the few-shot learning setting is similar to the performance in the one-shot learning setting. The overall ratio of matches and mismatches is 83:37, or around 69.2 per cent, however the composition is different, compared with the one-shot learning setting.

Table 6: Classification of sexist statements with one-shot learning

Actual classification	GPT-3 classification	
	Not sexist	Sexist
Not sexist	14	16
Sexist	7	23

Table 7: Classification of hate speech with one-shot learning

Actual classification	GPT-3 classification	
	Not hate speech	Hate speech
Not hate speech	35	25
Hate speech	12	48

Table 8: Performance of one-shot learning in hate speech classification

	Percent
Racism	
Accuracy	76.67
Precision	73.53
Recall	83.33
F1	78.12
Sexism	
Accuracy	61.67
Precision	58.97
Recall	76.67
F1	66.67
Overall	
Accuracy	69.17
Precision	65.75
Recall	80.00
F1	69.17

Table 9: Classification of racist statements with single-category few-shot learning

Actual classification	GPT-3 classification	
	Not racist	Racist
Not racist	18	12
Racist	7	23

Table 10: Classification of sexist statements with single-category few-shot learning

Actual classification	GPT-3 classification	
	Not sexist	Sexist
Not sexist	24	6
Sexist	12	18

Table 11: Classification of hate speech with single-category few-shot learning

Actual classification	GPT-3 classification	
	Not hate speech	Hate speech
Not hate speech	42	18
Hate speech	19	41

Table 12: Performance of single-category few-shot learning in hate speech classification

	Percent
Racism	
Accuracy	68.33
Precision	65.71
Recall	76.67
F1	70.77
Sexism	
Accuracy	70.00
Precision	75.00
Recall	60.00
F1	66.67
Overall	
Accuracy	69.17
Precision	69.49
Recall	68.33
F1	60.51

Few-shot learning – mixed category

The results of the mixed-category few-shot experiments are presented in Table 13. Among the ten sets of examples, Example Set 2 yields the best performance, as the model has the highest number of matches, 180, and the lowest number of mismatches, 60. In other words, the accuracy under the mixed category few-shot setting with Example Set 2 is 75 per cent. The example set that yields the worst results is Example Set 9, in which there are 137 matches and 103 mismatches. The accuracy rate with Example Set 9 is 57.1 per cent. The difference between Example Sets 2 and 9 suggests that, although the models are provided with same number of examples, the content of the examples also affects how the model makes inferences.

The unique generated answers are listed in Table 14. These are the response of GPT-3 that we obtain when we ask the model to classify statements, but do not provide examples that would serve to limit the responses. Under the mixed-category setting, the model is observed to generate answers that are out of scope. For instance, other than ‘sexist’, ‘racist’, and ‘neither’, we also see answers such as ‘transphobic’, ‘hypocritical’, ‘Islamophobic’, and ‘ableist’. In some cases, the model even classifies a text passage into more than one category, such as ‘sexist, racist’ and ‘sexist, misogynistic’. The full list contains eighty different answers instead of three.

Few-shot learning – mixed category with instruction

To reduce the chance of the model generating answers that are out of scope, a brief instruction is added to the prompt, specifying that the answers be: ‘sexist’, ‘racist’, or ‘nei-

Table 13: Classification of statements with mixed-category few-shot learning

Example Set	Result	Count	Percent
1	Match	154	64
1	Mismatch	86	36
2	Match	180	75
2	Mismatch	60	25
3	Match	167	70
3	Mismatch	73	30
4	Match	157	65
4	Mismatch	83	35
5	Match	157	65
5	Mismatch	83	35
6	Match	173	72
6	Mismatch	67	28
7	Match	144	60
7	Mismatch	96	40
8	Match	155	65
8	Mismatch	85	35
9	Match	137	57
9	Mismatch	103	43
10	Match	174	72
10	Mismatch	66	28

Table 14: Classifications generated by GPT-3 under mixed-category few-shot learning without instructions

racist | neither | homophobic | sexist | sexist, racist, | sexual assault | religious | sexual harassment | racist, sexist, | transphobic | hypocritical | I’m not a | I’m not brave | none of your business | I didn’t | I’m not offended | you’re not | sexual | I don’t play | I don’t care | I don’t know | victim blaming | it’s a question | irrelevant | I’m not crazy | creepy | romantic | I was taught to | I didn’t say | ignorant | I’m not sure | I agree | not true | not even wrong | YouTube doesn’t remove | socialist | conspiracy theorist | overpopulation | ableist | Islamophobic | conspiracy theory | hippie | sexist, rape ap | cliché | not even close to | nostalgic | I don’t think | not a question | hippy | terrorist | hate speech | he was a terrorist | preachy | not a valid argument | not sexist | brave | Islamophobe | because you’re a | troll | mental | never | pedophilia | mental illness is not | I’m sure you | I’m not going | not even close | grammar | rape apologist | cliché | vegan | sexist, misogynistic | I do | wrong | funny | question | when you were taught | opinion | emotional | conspiracy | nostalgia

Table 15: Classification of statements with mixed-category few-shot learning, with instruction

Example Set	Result	Count	Percent
1	Match	187	78
1	Mismatch	53	22
2	Match	184	77
2	Mismatch	56	23
3	Match	179	75
3	Mismatch	61	25
4	Match	185	77
4	Mismatch	55	23
5	Match	172	72
5	Mismatch	68	28
6	Match	170	71
6	Mismatch	70	29
7	Match	182	76
7	Mismatch	58	24
8	Match	174	72
8	Mismatch	66	28
9	Match	173	72
9	Mismatch	67	28
10	Match	179	75
10	Mismatch	61	25

ther’. The addition of an instruction successfully restricts the generated answers within the specified terms. The unique generated answers are: ‘racist’, ‘neither’, and ‘sexist’.

The results of the mixed-category, few-shot learning, with instruction, experiments are presented in Table 15. With the addition of an instruction in the prompt, the example set that yields the best result is Example Set 1 instead of Example Set 2. The addition of the instruction in the prompt increases the highest number of matches from 180 to 187; the highest accuracy rate is increased from 75 per cent to 78 per cent. In almost all cases the accuracy of the model increases (Figure 1). Across all the example sets, there are 1,785, or 74 per cent, matches and 615, or 26 per cent, mismatches.

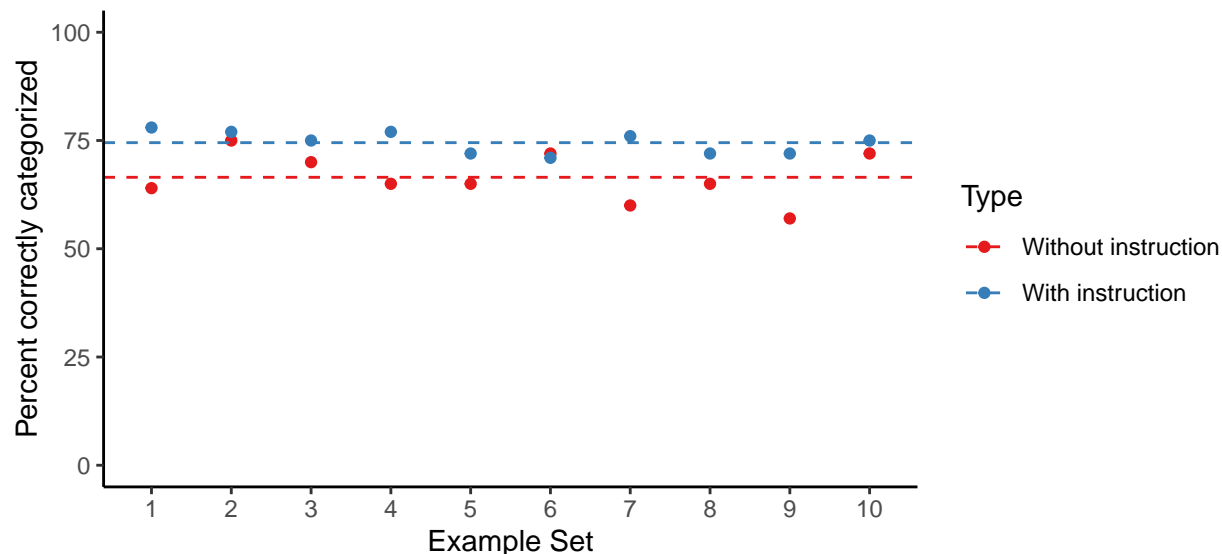


Figure 1: Comparing accuracy with and without an instruction

Discussion

In the zero-shot learning setting where the model is given no examples, its accuracy rate for identifying sexist and racist text is 48 per cent. In both the one-shot learning setting and the single-category few-shot learning setting, the accuracy increases to 69 per cent. It is likely that the model is not suitable for use in hate speech detection in either the zero-shot learning or one-shot learning settings, as the accuracy rates are between 48.3 per cent and 69 per cent. Davidson et al. (2017), using a different model and approach, similarly find ‘that almost 40 per cent of hate speech is misclassified’. And when Schick et al. (2021) use GPT-2 they find a similar ability to recognise sexually explicit content, however using an alternative model - Google’s T5 (Raffel et al., 2020) - they find better results.

In the mixed-category few-shot setting, different example sets yield different accuracy rates. Without an instruction included in the prompt, the accuracy rates range from 57 per cent (Example Set 9) to 75 per cent (Example Set 2). Adding an instruction improves the performance, and the highest accuracy rate increases from 75 per cent (Example Set 2) to 78 per cent (Example Set 1). Overall accuracy rates also increase in almost every example set after adding an instruction in the prompt. The large difference in accuracy rates between the example sets suggests that the content of the examples may matter. It may be that deliberately curated examples would enable better accuracy.

Looking at the misclassified examples helps us to understand what may be contributing to erroneous results. The following comment has a true label of ‘racist’, but the model classifies it as ‘sexist’. This comment uses racial identity terms such as ‘Asian’ and ‘White’, but its core argument is about discrimination on the basis of sex. The model did not have the necessary sensitivity to categorize the comment as sexist instead of racist.

‘The male gene is stronger, as it should be of course. Have an Asian male

breed with a White women and the child will look very asian. Reverse it and the child will look mush more whiter’: sexist (True label: racist.)

As another similar example, the following comment has a true label of ‘neither’, but the model classified it as ‘transphobic’.

‘Civil rights movement for transgendered people.’: transphobic (True label: neither.)

The comment is a plain statement without judgment; the inclusion of ‘transgendered people’ may be the reason for the false positive. If this were the case, then combined with the previous example, the model may have a tendency to generate false positive results when it encounters words that are identity terms, erroneously ignoring the context in which they are used. Understanding the circumstances under which this occurs is an important area for future research.

The following comment is an example of false negativity. This comment has a true label of ‘racist’ and was classified as ‘neither’ by the model. This is possibly because of the misspelling of the profane word ‘phuck’, where the implied word, which carries aggression, is obvious to a human reader.

‘phuck white rasist ass people’: neither (True label: racist.)

If this were the case, then it also points to a weakness of the use of this type of model for detecting hate speech in human-generated content. In particular, the way that humans use words changes quickly, especially sexist and racist language online. Large language models take a long time to train and deploy. This means that it may be that unless considerable context is provided in the prompt, that the model would not identify words that have become sexist and racist through their usage in the time since the model was trained.

In conclusion, with proper settings such as the inclusion of instruction and curated examples, large natural language models such as GPT-3 are capable of identifying sexist and racist text at a similar level of specificity to other methods. However, it is possible that if a user intentionally misspells profane words, the models may be less likely to identify such content as problematic. This possibility deserves further investigation due to the tendency for language to change quickly. Furthermore, the models might classify texts that contain identity terms as problematic, as they are often associated with harmful text when being targeted. Various settings should be further explored to tackle these obstacles. To shed light on why a text is misclassified, it might be appropriate to prompt GPT-3 to deliver an accompanying explanation for the decision, or to consider a larger dataset of sexist and racist content. The extent to which the identification and classification of hate speech can be explained is especially of interest for future work and is something that could be implemented on the publicly available GPT-2 from OpenAI or T5 from Google.

Appendix

Table 16: The ten example sets for the mixed-category, few-shot learning, experiments

Set	Example 1 (sexist)	Example 2 (racist)	Example 3 (not hate speech)
1	You should know women's sports are a joke	why does that african american monkey hold gun like that?	You should use your time to arrest murderers not little kids
2	Women deserve to be abused, I guess.	Why do negroes tear down their own neighborhoods when a black thug gets killed by law enforcement? Must be because of low IQ.	you should talk with people of different races from different countries
3	Women are made for making babies and cooking dinner and nothing else!!!	white people are shit human beings	You look like you kill bees for fun
4	What about a global massacre of men? Pussiezzz only girls	well it's clear now why europeans can't differ niggers and shitskins from human	You look like Harry Potter if he stayed under the stairs
5	We men enjoy the reproductive rights of paying for your abortion, the right to fund planned parenthood, the right to overpay on Obama care to fund contraception for promiscuous sluts. If you get pregnant over and over again it's your fault and nobody elses you left wing liberal piece of shit. Go be a whore. Keep spreading disease. Half of all black women have herpes. And we have the right to pay for it. You are outrageous.	We whites have higher IQ's than blacks, this is because IQ measures intelligence and we are the more evolved and superior race with a larger brain.	You look kinda dead inside, you parents must not haved loved you.
6	The men are turtles...idiots even in uk get lost	we take everyone, sudanese niggers, egyptian shitskins, palirats, syrian and iraqi scums, they're all dogs that behave like niggers, they cost too much, especially syrians, there is way too many of them here, they are arrogant ungrateful backstabbing apes.	You have the grace and composure of a rabid hyena
7	That girl still need a good ass whooping, stupid ass bitch!!	U SUCK SHIT YOU BLACK FUCK	You deserve way more subscribers . The shit you say is fucking gold.
8	stone throwing must be legal when against women	This sick, disgusting, white terrorist/fascist killed more than 50 innocent people.	You can see the country is fertile and green.ă just need the right governance
9	She is heavily relying on him to turn the other cheek...tough talking demon infested woman.	This is considered WHITE in Sweden. No wonder they accepted all those black and muslim immigrants. They are literally color blind.	you are a great brother to abbie, that family is changing the world!
10	Really, I hope she get raped in the future. Raping women should be legal I think	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 adressñññÜñóØññèñüñññØñññðñÇ	You are a brave man.....for letting them keep the comments enabled

References

- Anagnostou, A., Mollas, I., and Tsoumakas, G. (2018). Hatebusters: A Web Application for Actively Reporting YouTube Hate Speech. In *IJCAI*, pages 5796–5798.
- 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*.
- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language model. *Journal of Machine Learning Research*, 3(Feb):1137–1155.
- 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*.
- Criminal Code (1985). Government of Canada. As viewed 19 March 2021, available at: <https://laws-lois.justice.gc.ca/eng/acts/c-46/section-319.html>.
- Davidson, T., Bhattacharya, D., and Weber, I. (2019). Racial bias in hate speech and abusive language detection datasets. In *Proceedings of the Third Workshop on Abusive Language Online*, pages 25–35.
- Davidson, T., Warmesley, D., Macy, M., and Weber, I. (2017). Automated hate speech detection and the problem of offensive language. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 11.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- 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*.
- Hovy, D. and Spruit, S. L. (2016). The social impact of natural language processing. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 591–598.
- Kennedy, B., Atari, M., Davani, A. M., Yeh, L., Omrani, A., Kim, Y., Coombs, K., Havaladar, S., Portillo-Wightman, G., Gonzalez, E., et al. (2018). The gab hate corpus: A collection of 27k posts annotated for hate speech.
- Lin, S., Hilton, J., and Evans, O. (2021). Truthfulqa: Measuring how models mimic human falsehoods.
- McGuffie, K. and Newhouse, A. (2020). The radicalization risks of GPT-3 and advanced neural language models. *arXiv preprint arXiv:2009.06807*.

- Mollas, I., Chrysopoulou, Z., Karlos, S., and Tsoumakas, G. (2020). ETHOS: An Online Hate Speech Detection Dataset.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer.
- Rosenfeld, R. (2000). Two decades of statistical language modeling: Where do we go from here? *Proceedings of the IEEE*, 88(8):1270–1278.
- Schick, T., Udupa, S., and Schütze, H. (2021). Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in nlp.
- Schmidt, A. and Wiegand, M. (2017). A survey on hate speech detection using natural language processing. In *Proceedings of the fifth international workshop on natural language processing for social media*, pages 1–10.
- Srba, I., Lenzini, G., Pikuliak, M., and Pecar, S. (2021). Addressing hate speech with data science: An overview from computer science perspective. *Hate Speech - Multidisziplinäre Analysen und Handlungsoptionen*, page 317336.
- Turian, J., Ratinov, L., and Bengio, Y. (2010). Word representations: A simple and general method for semi-supervised learning. In *Proceedings of the 48th annual meeting of the association for computational linguistics*, pages 384–394.
- Twitter (2021). Hateful conduct policy. As viewed 19 March 2021, available at: <https://help.twitter.com/en/rules-and-policies/hateful-conduct-policy>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Waseem, Z. and Hovy, D. (2016). Hateful symbols or hateful people? Predictive features for hate speech detection on Twitter. In *Proceedings of the NAACL student research workshop*, pages 88–93.