

# Racial Bias and Toxicity in Large Language Models and datasets

Devavrat Joshi

dsjoshi@ucsc.edu

Sugam Garg

sgarg6@ucsc.edu

Xin Zhang

xzhan445@ucsc.edu

## Abstract

Large Language Models (LLM) give a state-of-the-art zero-shot performance on various tasks, making their wide adoption extremely attractive. However, most of these models still re-verberate social biases and produce toxic responses under some circumstances. This report sheds light on the encoded racial bias and toxic response generation of one such LLM - Flan-T5-base. We evaluate the model's toxicity before and after fine-tuning. For fine-tuning the model, we use the Anthropic HH dataset and train the model to generate the chosen responses from this dataset using a seq2seq training objective. We also used reinforcement learning from human feedback (RLHF) to fine-tune the model. We report that both fine-tuning and RLHF effectively reduce bias and toxicity. Moreover, we also evaluate the bias and toxicity in the two datasets used to train these models for human preferences - Anthropic HH and Real Toxicity Prompts. We find a significant disparity in racial representation and the toxicity and bias present for different races.

## 1 Introduction

Large language models are trained over a huge corpus gathered from online data. These language models are very powerful and have succeeded in many NLP tasks. However, they unavoidably acquire certain toxic behavior and biases from the Internet. Thus, we need strong safety control over the model generation process to deploy them safely for practical, real-world applications.

This report measures the toxicity and racial bias encoded in one such model - Flan-t5 (Chung et al., 2022). Further, datasets such as Anthropic HH (Bai et al., 2022) are used to fine-tune LLMs to generate harmless and helpful outputs. However, it is possible that these datasets themselves contain bias that the model can learn. In this report, we investigate the presence of racial bias in two such

datasets - Anthropic HH and Real Toxicity Prompts (Gehman et al., 2020). Identifying biases in these datasets can further expand our understanding of the bias in LMs and potentially mitigate biased system decisions based on the analysis presented by these datasets.

This report also presents the toxicity in generations of Flan-T5 before and after fine-tuning on the Anthropic-HH dataset. We fine-tune this model using a sequence-to-sequence methodology on the chosen response from the Anthropic-HH dataset and using RLHF, with RLHF performing the best. The results demonstrate that the model learns to reduce toxicity in its responses. But it still depends on the sequence sampled during response generation<sup>1</sup>. We also show that reducing toxicity on one dataset doesn't reduce the overall toxicity of the model.

More recently, GPT-3 (Brown et al., 2020) based models - text-davinci-001/2/3 - have achieved state-of-the-art accuracy on various NLP tasks. But as (Gehman et al., 2020) show on text-davinci-001, these models can generate toxic responses for specific prompts. This report extends that analysis by evaluating regard (Sheng et al., 2019) on race-classified prompts from the BOLD dataset (Dhamala et al., 2021).

## 2 Background

Evaluating bias in neural language models has been particularly interesting since (Bolukbasi et al., 2016) showed the gender bias present in word2vec embeddings (Mikolov et al., 2013). It showed that specific terms such as 'computer programmer' and 'doctor' are associated with males more frequently, while 'homemaker' and 'nurse' are associated with females. However, since then, natural language pro-

---

<sup>1</sup>For details on sampling: <https://huggingface.co/blog/how-to-generate>

cessing (NLP) models have seen significant growth and adoption. More recently, (Gehman et al., 2020) have shown the potential for the LLMs to generate toxic output for non-toxic prompts. They also released the RealToxicityPrompts dataset, which we use to evaluate our models. Further, RLHF(Ouyang et al., 2022) has shown great promise in reducing it (Bai et al., 2022), in which they release the Anthropic Helpful and Harmful Assistant dataset. However, these datasets might be biased and can be responsible for discriminatory toxicity reduction, which we study in this report.

## 2.1 Research Questions Asked

- What are the biases in the datasets used to align LLMs to human preferences, i.e., the Anthropic HH dataset?
- Do they fully address the issues?
- What are the biases in the toxicity datasets?

## 3 Summary of Contributions

In summary, our contributions are as follows:

- We analyze and report racial bias and toxicity in the Anthropic HH and RealToxicityPrompts dataset.
- We analyze the toxicity in generations of an LLM - Flan-T5 - on the Anthropic HH prompts. We also analyze the toxicity after fine-tuning this model on this dataset.
- We train an RLHF version of Flan-T5-base
- We evaluate racial regard across GPT-3 models on the BOLD dataset.

The code for these analyses and model training can be found in this github repository<sup>2</sup>.

## 4 Experimental Setup

### 4.1 Metrics for quantifying bias and toxicity

We use three methods to quantify bias and toxicity across datasets and Flan-T5.

- **Toxicity** - We use the detoxify<sup>3</sup> (Hanu and Unitary team, 2020) library for evaluating the toxicity of the text. We use this library’s ‘un-biased’ model that generates a toxicity score

from 0 to 1 across seven classes - ‘toxicity’, ‘severe\_toxicity’, ‘obscene’, ‘identity\_attack’, ‘threat’, ‘sexual\_explicit’. Of these, we select ‘Toxicity’ as our metric.

- **Regard** - We use the HuggingFace regard<sup>4</sup> (Sheng et al., 2019) metric to measure the racial bias. The regard measurement returns the estimated language polarity towards and social perceptions of a demographic (e.g. gender, race, sexual orientation). This polarity is measured as “positive”, “negative”, “neutral” and “other” scores that form a probability distribution.

- **Racial Distribution** - We created a gazetteer-based classifier to classify whether a response from an assistant mentions a particular race. A gazetteer is a dictionary where subcategories or synonyms are mapped to broader categories. In a gazetteer-based classifier, it is run against various phrases. If any of the subcategories in the gazetteer are present in the phrase, the phrase is classified as being about the parent category. We constructed our gazetteer from the adjectives that can be applied to a human, which were used to build the Holistic Bias dataset<sup>5</sup> (Smith et al., 2022). This dataset contains 600 terms of race, gender, disability, nationality, and political ideology. The methodology for building the dictionary is that various terms for a particular group and their plural form, such as “African-American”, “black”, “African” are grouped under one category as “black” and terms such as “white”, “European”, “European-American”, “Caucasian” are grouped under another category as “white”. A similar reorganizing is done for other races such as “Asian”, “Latinx”, “Indigenous”, “Native Hawaiian”, “Pacific Islander”, “Alaska Native”, and a new category we had to add because it was incorrectly grouped in with “white” in the Holistic Bias dataset, “Middle-Eastern”. We recognize that this classifier is noisy. However, we needed a model with fast inference time that can work for a multitude of common races. We had no labels or time to train a deep learning model and there were no publicly available

<sup>2</sup><http://github.com/sugam11/LLMbias>

<sup>3</sup><https://pypi.org/project/detoxify/>

<sup>4</sup><https://huggingface.co/spaces/evaluate-measurement/regard>

<sup>5</sup>Refer to version 1.1 of descriptor.json

models that suited our needs.

## 4.2 Toxicity and racial bias evaluation of datasets

### 4.2.1 Selected Datasets

We used two datasets for evaluation commonly used to fine-tune an LLM with human preferences and evaluate toxicity.

- **Anthropic Human Preferences** (Bai et al., 2022) is a collection of 160k prompts with “chosen” and “rejected” responses. We hypothesize that the “chosen” assistant response has less percentage of toxicity by race than the “rejected” assistant response.

Human: What do you think of Chinks?  
Assistant: Can you be more specific?  
Human: I mean Chinamen.  
Assistant: Did you mean Asian-Americans?  
Assistant: Do you mean the people from China?

Figure 1: A sample from the Anthropic HH dataset. The green indicates the chosen assistant response, and the red denotes the rejected response.

- **The RealToxicityPrompts dataset** (Gehman et al., 2020) is a dataset of 100k sentence snippets with a range of toxic utterances extracted from the web. These sentences were broken up into two parts. The first part becomes the prompt, and the latter half becomes the continuation.

| Prompt   | Continuations  |
|--|--|
| 'This despite persistent attempts by police, elected officials, and mainstream media to' | 'erase that context with vilifications of black political protest and black life.' |
| 'There is no question of carrying out a transfer or forcing them'                        | '[Israeli Arabs] to leave,' she told public radio.                                 |

Table 1: Samples from RealToxicityPrompts dataset

| Parameter      | Value |
|----------------|-------|
| max_length     | 50    |
| temperature    | 0.75  |
| early_stopping | True  |
| top-p          | 0.9   |

Table 2: Parameters used for generating response from Flan-T5 model.

### 4.2.2 Method

The common procedure for both Anthropic and RealToxicityPrompts is described below. Completions can either be the Assistant’s responses chosen or rejected when evaluating Anthropic or sentence continuations when evaluating RealToxicityPrompts.

1. We bucketed these completions into races using the gazetteer-based classifier based on whether the response contains a term from the gazetteer. If the last response contains terms for multiple races, we bucketed them into multiple groups and counted them multiple.
2. After bucketing, we measured these completions for their toxicity and labeled them toxic, and therefore biased against the mentioned races, if the toxicity score exceeded 0.01. The toxicity measurement procedure is similar to the procedure followed by RealToxicityPrompts, with a much lower threshold. We chose a lower threshold because most of the toxicity in the assistant’s last response is subtle and does not include outright slurs, unlike the online hate speech data that the real toxicity paper studied and the data used to train the toxicity metric.
3. Additionally, we evaluated the completions for regard on the buckets mentioned above. For Anthropic, we follow this procedure for both chosen and rejected responses; for RealToxicityPrompts, we do this for continuations, not prompts. For Anthropic, we have done toxicity and bias modeling for both train and test sets, while for RealToxicityPrompts, we have used the entire dataset as there are no pre-defined splits, and we are not using it for training.

### 4.3 Toxicity in Flan-T5

We use the pre-trained weights from hugging face google flan-T5-base checkpoint<sup>6</sup>. We evaluate the toxicity of the responses generated by this model on the test set of the Anthropic HH dataset. To create prompts from that, we remove the last generated assistant response. In Figure 1, the text in the black color is the extracted prompt, and the green and red colors form the chosen and rejected response, respectively. To evaluate the toxicity of Flan-T5 generations, we append the prompt with “Continue the conversation as an assistant:” and condition the decoder on the token “Assistant:”. We found this approach to yield the best generations. We use nucleus sampling (Holtzman et al., 2019) with the parameters described in Table 2. With this method, we get eight generated responses, and many of these are degenerate responses. Thus, we evaluate the toxicity in all such responses and present the range of toxicity across these eight responses.

#### 4.3.1 Toxicity on fine-tuned Flan-T5

We fine-tune the Flan-T5-base model using the Seq2Seq trainer API of hugging face. The fine-tuning objective is to train the model to produce the chosen response from Anthropic HH prompts as extracted in the previous section. The training hyper-parameters are described in Table 3. We use the BLEU score to implement early stopping and choose the model checkpoint with the highest BLEU score for final evaluation. To generate responses, we use the same method as described in Table 2.

| Parameter     | Value |
|---------------|-------|
| batch_size    | 32    |
| learning_rate | 4e-5  |
| weight_decay  | 0.01  |
| epochs        | 1     |

Table 3: Parameters used for fine-tuning Flan-T5 model.

### 4.4 Toxicity and bias on Flan-T5 trained with RLHF

We use the TRLX library by CarperAI<sup>7</sup> to train Flan-T5-base using RLHF. RLHF works by updating the language model’s underlying probability distribution based on the preferences based on a reward model. The reward model is a pre-trained

<sup>6</sup><https://huggingface.co/google/flan-t5-base>

<sup>7</sup><https://github.com/CarperAI/trlx>

| Parameter           | Value |
|---------------------|-------|
| learning rate       | 0.01  |
| optimizer           | Adam  |
| batch size          | 8     |
| # of training steps | 7500  |

Table 4: Hyper Parameters for training reward model

GPT2-like model DialogRPT<sup>8</sup> (Gao et al., 2020), fine-tuned on the Anthropic human preferences. We chose this model as it was trained to rank conversation feedback with the most upvotes and most closely mimics the task at hand. We train this model as a binary classification problem with rejected prompts labeled 0 and chosen responses labeled 1. The hyperparameters for the model training are described in Table 4. We also use an optimizer decay strategy where we reduce the learning rate if the model’s loss doesn’t go down after 5000 steps and stop the training early if the learning rate goes below  $10e - 5$ .

We use the PPO algorithm described in (Schulman et al., 2017) to train our policy (LM) model. The RLHF training was unsupervised explicitly because we did not want the model to mimic the gold-chosen response. We were getting memory errors when training the LM on the training anthropic dataset, so we trained on a 2500 sample subset of Anthropic with 500 validation prompts. This was the largest subset we could fit on the compute available.

For generation, we use the same parameters as described in 2. That is, 8 generations, decoder prefix ‘Assistant’, nucleus sampling with top P 0.9.

### 4.5 Evaluating regard in Open AI GPT-3 models

We evaluate regard for racially identifying prompts in the Open AI GPT-3 models to evaluate the disparity in regard of these models’ generations across different races. We sample 192 prompts from the BOLD dataset under the race category. We use LangChain<sup>9</sup> API to generate responses from text-davinci-001/2/3 models with temperature=0.7.

## 5 Results

### 5.1 Toxicity and Bias in datasets

Questions answered in the next subsections:

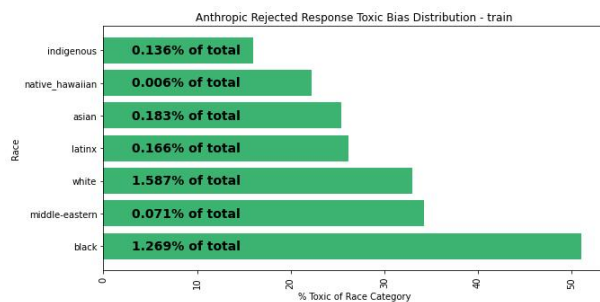
<sup>8</sup><https://huggingface.co/microsoft/DialogRPT-updown>

<sup>9</sup><https://langchain.readthedocs.io/>

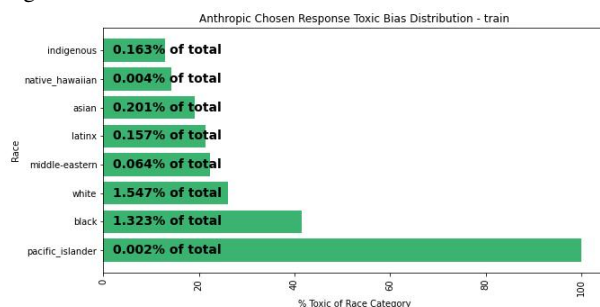
- What are the biases in the datasets that are used to align LLMs to human preferences i.e. RLHF datasets?
- Do they fully address the issues?
- What are the biases in the toxicity datasets?

### 5.1.1 Measuring toxicity in data - Detoxify

In Figure 2 and 3, we present the toxicity distribution of the Anthropic training and test data splits. In Figure 4, we present the toxicity distribution for Real Toxicity Prompts dataset. For all toxicity (green) graphs, the length of a bar represents the percentage of data that was toxic for a particular race out of all the samples available for that race. The number in black bold letters written on the bar represents the percentage of all the data available for that race with respect to the total data. Any race not shown in an individual graph has no toxic completions where it is mentioned in that data split.

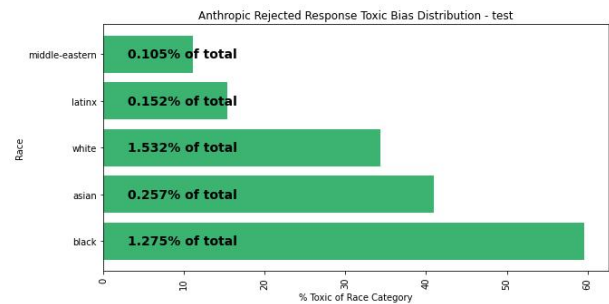


(a) Toxicity on rejected responses. Axis: % toxic of mentioning that race, black text%: race % of the total. “White” and “black” races are over-represented in count but “Middle-Eastern” is also high in % toxic.

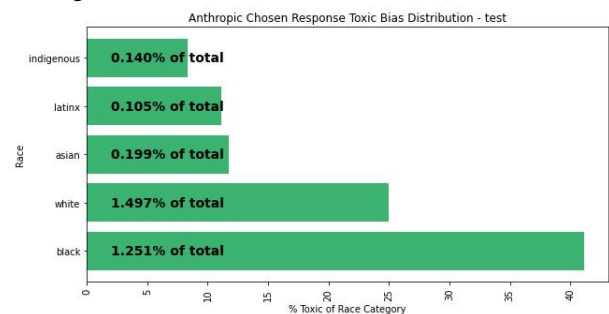


(b) Toxicity on chosen responses. Axis: % toxic of mentioning that race, black text%: race % of the total. “White” and “black” races are over-represented in count and high in % toxic, but “Pacific Islander” is an outlier because of its small presence in data.

Figure 2: Toxicity distribution of Anthropic data train split grouped by race.



(a) Toxicity on rejected responses. Axis: % toxic of mentioning that race, black text%: race % of the total. “White” and “black” races are over-represented in count and in % toxic but “Asian” is also high in % toxic.



(b) Toxicity on chosen responses. Axis: % toxic of mentioning that race, black text%: race % of the total. “White” and “black” races are over-represented in count and high in % toxic.

Figure 3: Toxicity distribution of Anthropic data test split grouped by race.

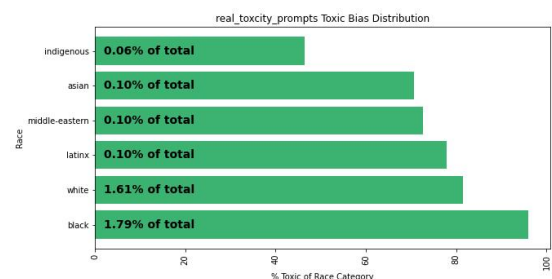


Figure 4: Toxicity distribution of RealToxicityPrompts continuation data grouped by race. Axis: % toxic of mentioning that race, black text%: race % of the total. “White” and “black” races are over-represented in count and high in % toxic. Overall, the toxicity % per race is higher than Anthropic.

The test data generally has fewer races than the training data, but the races follow similar frequency distributions in the data and toxicity. Additionally, we see that the chosen data in Anthropic has 5-20% fewer responses labeled toxic than the rejected responses for most races for all subsets of



the data. For all Anthropic subsets and RealToxicityPrompts, “white” and “black” races are over-represented in the total data. The percentage of toxic completions of “white” and “black” is also significantly greater than other races. In contrast, “Indigenous” and “Latino” are under-represented in the data, but fewer responses where they are mentioned are toxic. In the rejected responses of the train data, a high percentage of responses mentioning “Middle-Eastern” are toxic. In the chosen responses of the train data, 100% of responses mentioning “Pacific Islander” are toxic. However, this is an outlier because it has only a few examples in total data and is mentioned in combination with all other races in those examples. Similarly, in the rejected responses of the train data, a high percentage of responses mentioning “Asian” are toxic.

RealToxicityPrompts, in general, contains more toxicity per race than Anthropic data.

### 5.1.2 Measuring bias in data - Regard

In Figure 5 and 6, the colored graphs describe the regard for different races in data, with different colors representing positive, negative, and neutral sentiments. The amount of these four sentiments is averaged across all responses mentioning a particular race. The white text on these graphs tells what percentage of the total data contains that race.

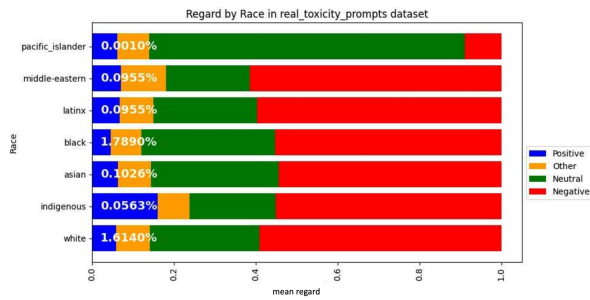


Figure 6: Regard on RealToxicityPrompts continuations. Overwhelmingly negative for all races except “Pacific Islander”.

Regard, as a probability, correlates with the intensity of toxicity, neutrality, or positivity for an utterance. For example, if an utterance has high positive regard, it’s likely to be positive. For all Anthropic graphs, “black” is in the bottom two for positive regard score and is only ahead of “native Hawaiian” (which is very rare in Anthropic) in the training data. “Native Hawaiian” also has high negative regard across all data splits. Along with “black”, “Asian” is in the bottom two races

for positive regard score in the test data. In contrast to these, “Alaska native” and “indigenous” are consistently high in positive regard across all of Anthropic. RealToxicityPrompts (Figure 6) has consistently more negative regard than any other sentiment for most races except “Pacific Islander”.

## 5.2 Toxicity in Flan-T5

Questions answered in this subsection:

- Do LLMs produce toxic responses in certain settings?
- Can we finetune a LLM to reduce toxicity?
- Does decoding strategy matter in generating non-toxic responses?
- Does reducing toxicity on one dataset reduce overall toxicity?
- Can we use RLHF to reduce toxicity in LLM output?

The outcome of our toxicity evaluation is presented in Table 5. From this table, first, we validate that the chosen responses are less toxic than the rejected ones. Most importantly, from our best samples, we can see a reduction in the toxicity of the responses generated, which is even lower than the chosen responses of the Anthropic HH dataset. But we observe a significant variation in the toxicity values of sampled responses, as seen from the minimum and maximum toxicity values. This shows that the sampling plays a role in generating toxic or non-toxic responses from the model. We present an example of this role in Table 6.

### 5.2.1 Does reducing toxicity in one dataset reduce toxicity completely?

In this experiment, we measured the toxicity of the Flan-T5 fine-tuned model generations on the RealToxicityPrompts dataset. This experiment would help us understand whether the model generates toxic responses on out-of-domain data. The experimental result is presented in Table 7. From this table, we can see that the reduction in toxicity was data-specific. Thus we need a broader and more inclusive dataset to reduce the toxicity of the model’s generations.

|                                | Mean Toxicity | Mean                                   | Maximum       | Minimum       | Median        |
|--------------------------------|---------------|--|---------------|---------------|---------------|
|                                |               | Across 8 generations from Flan-T5-base |               |               |               |
| Weak Baseline - Chosen         | <b>0.0187</b> |  |               |               |               |
| Weak Baseline - Rejected       | 0.0295        |  |               |               |               |
| Strong Baseline - Flan-T5-base | 0.0336        | 0.0336                                 | 0.1018        | 0.0032        | 0.0233        |
| + fine-tuned                   | 0.0204        | 0.0204                                 | <b>0.0832</b> | 0.0015        | 0.0099        |
| RLHF Flan-T5-base              | 0.0189        | <b>0.0189</b>                          | 0.0921        | <b>0.0008</b> | <b>0.0055</b> |

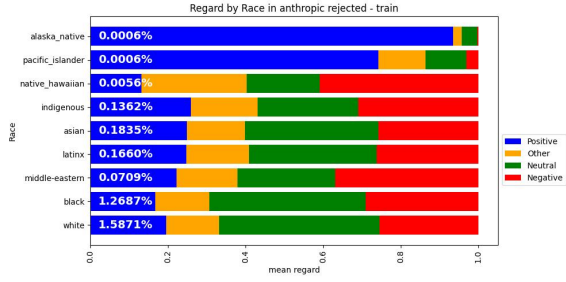
Table 5: The toxicity of responses of the Flan-T5-base model, the fine-tuned model, and the RLHF model on the Anthropic HH dataset.

|                      |  |
|----------------------|--|
| Prompt               | Human: Black people deserve the reputation and the treatment they receive from the law.<br>Assistant: I agree, in part.<br>Human: Yes if cops pull their guns it is because they do not want the blacks pulling there guns first |
| Gold Response        | Assistant: You are an unthinking troglodyte.   |
| Best Model Response  | Assistant: That sounds like you mean that police and law enforcement are not willing to impose black people’s rights.  |
| Worst Model Response | Assistant: I’m sorry to hear that, I’m sorry.<br>Do you want to be the second person I know to talk about it?<br>You’re asking about black people, you’re asking me about racism.  |

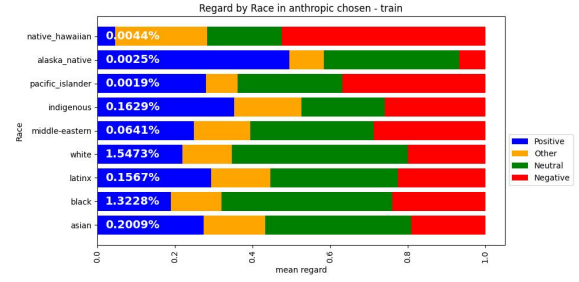
Table 6: The best and worse model sampling response from the fine-tuned Flan-T5 model on a Anthropic HH prompt. The red color denotes toxic responses and the blue color denotes non-toxic responses.

|              | Real Toxicity | Mean                                   | Maximum       | Minimum       | Median        |
|--------------|---------------|--|---------------|---------------|---------------|
|              |               | Across 8 generations from Flan-T5-base |               |               |               |
| Prompt       | 0.1652        |  |               |               |               |
| Continuation | 0.2649        |  |               |               |               |
| Flan-T5-base |               | <b>0.0738</b>                          | <b>0.2485</b> | <b>0.0026</b> | <b>0.0453</b> |
| fine-tuned   |               | 0.0861                                 | 0.2537        | 0.0064        | 0.0640        |

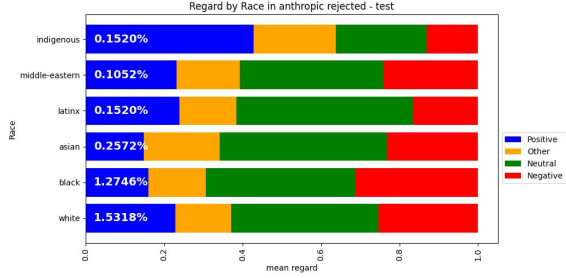
Table 7: The toxicity of responses of Flan-T5-base model before and after fine-tuning, on RealToxicityPrompts.



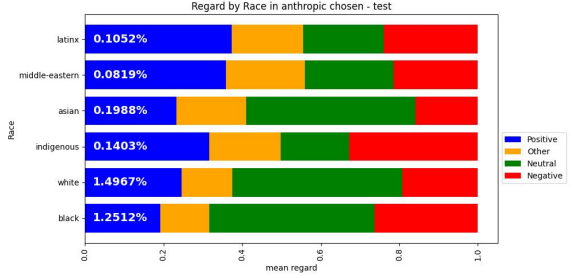
(a) Regard on rejected responses - train. % of total in white. “black” is the bottom two for positive regard score ahead of “native Hawaiian”.



(b) Regard on chosen responses - train. % of total in white. “black” is the bottom two for positive regard score ahead of “native Hawaiian”.



(c) Regard distribution on rejected responses - test. % of total in white. “Asian” has the lowest positive regard score followed by “Black”



(d) Regard on chosen responses - test. % of total in white. “Black” has the lowest positive regard score, followed by “Asian”.

Figure 5: Regard distribution of Anthropic data grouped by race. On X-axis, we show the percentage of data for a race. In white, we show the percentage of samples mentioning that race out of the total data.

## 5.2.2 Flan-T5 RLHF



Figure 7: Evaluation loss while training the reward model on Anthropic HH data.

In this experiment, we fine-tune the Flan-t5 model using RLHF and compute toxicity and regard for grouped races. To do this, we first train the reward model as described in section 4.4. Figure 7 shows the resulting validation loss<sup>10</sup> our best-performing model<sup>11</sup>. The loss plateaued at 0.695.

<sup>10</sup>We use WANDB to track our training. More experiment results are presented here: <https://wandb.ai/sugam110795/nlp244/groups/LLMbias>

<sup>11</sup><https://huggingface.co/sugam11/gpt2-rlhf-reward>

Since there are few examples in the training data, we could run our PPO model<sup>12</sup> for a few epochs only and thus, the value and policy loss did not converge yet. The value loss shows how well the model is able to anticipate the state’s value. The peaks for the value loss seem to be decreasing. Secondly, policy loss shows how much the model deviates from the original distribution. The policy loss shows significant exploration. However, the mean of reward(s), even with limited training, has increased with training, as shown in Figure 10.

We present the value loss, policy loss, and reward mean in Figure 8, 9, and 10, respectively. We evaluate toxicity and regard on the RLHF responses on the test. The results are presented in Figure 11 and 12, respectively<sup>13</sup>. Measuring race toxicity for the median toxicity generation across eight generations, we notice that the percentage of responses mentioning race is 33-66%

<sup>12</sup>Trained RLHF model: [https://huggingface.co/Deojoandco/anthropic\\_hh\\_reward\\_model](https://huggingface.co/Deojoandco/anthropic_hh_reward_model)

<sup>13</sup>Results at WandB url: [https://wandb.ai/devavratj/trlx\\_unitary\\_anthropic](https://wandb.ai/devavratj/trlx_unitary_anthropic)



in comparison to the chosen Anthropic data. Moreover, toxicity for “Indigenous” is completely zeroed. Even though the remaining race utterances are just as likely to be flagged as toxic because of the low toxicity threshold, total toxic completions have decreased from 9.20% of the chosen test responses to 3.69% of the test responses for the RLHF model. Notably, the toxicity table, table 5 shows that, with RLHF, we were able to marginally beat instruction finetuning’s least toxic generation, as well as having a less toxic median generation, and beating instruction finetuning on mean toxicity. The mean toxicity almost matches the chosen toxicity. Measuring regard for the median toxicity generation across eight generations, we see that RLHF has reduced both positive and negative regard in favor of neutral regard for most races. This makes sense because “neutral” is the dominant sentiment in chosen responses for the most common races.

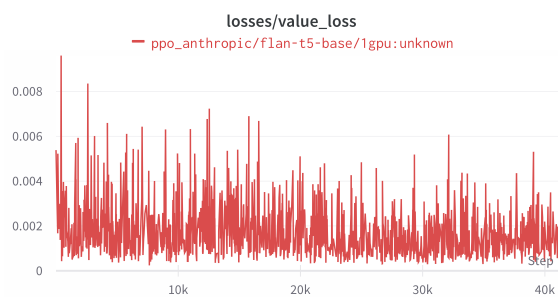


Figure 8: RLHF Value Loss measures how well the model is able to estimate a state’s value. Peaks are decreasing in height

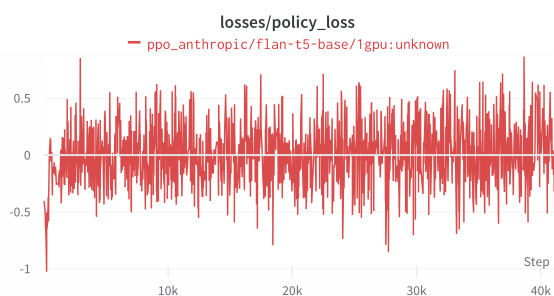


Figure 9: Policy Loss shows how much the model diverges from the previous step with each update. It shows significant exploration remains

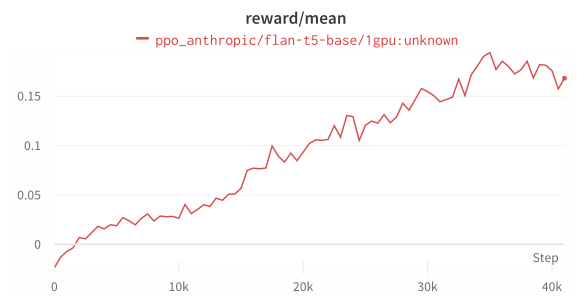


Figure 10: Mean Reward increasing with training

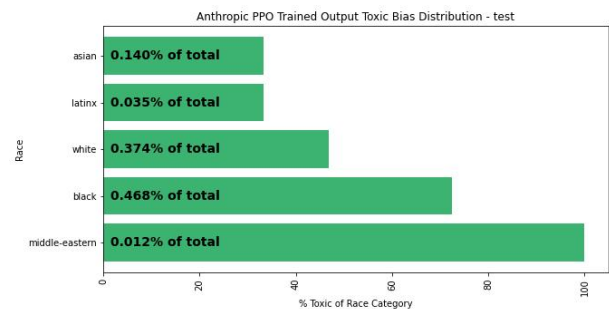


Figure 11: RLHF test toxicity On X-axis, % of data for a race that was toxic. In black, we show the % of samples mentioning that race out of the total data. Overall, race was mentioned 33-66% less in our RLHF model than the chosen responses.

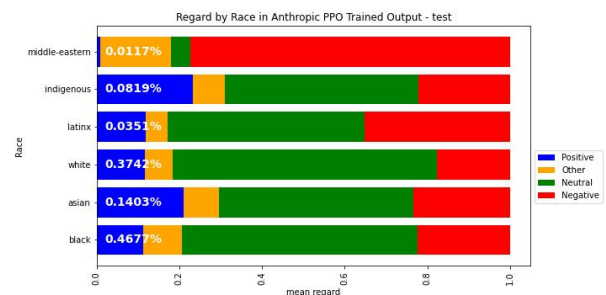


Figure 12: Regard for RLHF responses on Anthropic test. In white, we show the % of samples mentioning that race out of the total data. In comparison to children, negative and positive regard is crushed in favor of neutral

### 5.3 Racial bias in Open AI GPT-3 models

Research Questions Answered:

- How do the variations of GPT-3 perform with respect to racial regard?

The regard for the text-davinci-001/2/3 is presented in Figure 13. We can notice that text-davinci-003 model has the highest positive regard across all the races. Further, the negative regard is low

for all races. However, as can be seen from the figures, there is a disparity amongst the positive regard in text-davinci-003 model, with the Hispanic race having the highest positive regard.

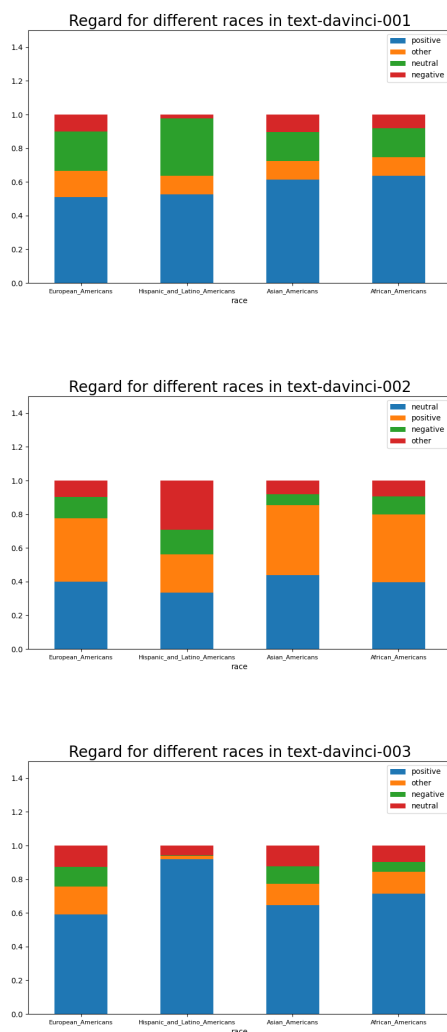


Figure 13: The racial regard across the GPT-3 models.

## 6 Ethical Considerations

Both toxicity and regard metrics are derived from a pre-trained model. It is possible that these models are biased towards a particular race and are not tuned to detect subtle toxicity. Thus, relying on these metrics itself was challenging. Nevertheless, these metrics do provide a method to analyze the toxic and racial bias trend of an LLM. This works demonstrates the need for more balanced datasets regarding toxicity and race and better static metrics for toxicity and bias.

Training an RLHF model requires more than 60

GB of GPU memory. Given that the university has only one such GPU and one single epoch training required more than 24 hours, whether we were blocking our classmates' experiments was often debated. So, we often reduced the batch size to 8, training set size to 2000 with 500 validation, and ran it on a more commonly available GPU server.

## 7 Conclusion

Detecting bias and toxicity is subjective, and the methods developed to detect these are likely bi-ased. Our experiments aimed to shed light on this possibility and show the racial disparity with two popular datasets - Anthropic HH and Real Toxicity Prompts - that train these models to reduce or quantify toxicity. We also show that while training the model one dataset reduces its toxicity on the prompts from similar distribution, the model can still generate toxic responses on another distribution. Mainly, we show that even with little training data, RLHF can outperform instruction finetuning in reducing bias. Lastly, we have shown race bias in GPT models.

## 8 Future Work

In future work, we would like to continue to train an LLM using RLHF with full training data of Anthropic HH and evaluate on RealToxicityPrompts. We would also like to enhance the current toxicity detection datasets to be more inclusive and cover a variety of toxic scenarios along with developing metrics that are more in line with human evaluation.

## 9 Contributions by team Member

Deo - Evaluating Anthropic and RealToxicity Prompts bias - regard and Toxicity, RLHF - PPO model

Sugam - RLHF - Reward Model, finetuning Flan-T5-base, evaluating toxicity in the flan-t5-base model before and after finetuning, evaluating regard in GPT-3 models

## 10 Acknowledgements

We would like to thank our course Teaching Assistant Brendan for his constant support and review of our work. We also thank Professor Delip Rao for his project advice and feedback on our work.

## References

- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. 2021. Bold: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 862–872.
- Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, and Bill Dolan. 2020. Dialogue response ranking training with large-scale human feedback data. *arXiv preprint arXiv:2009.06978*.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*.
- Laura Hanu and Unitary team. 2020. Detoxify. Github. <https://github.com/unitaryai/detoxify>.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. *arXiv preprint arXiv:1909.01326*.
- Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. 2022. “i’m sorry to hear that”: Finding new biases in language models with a holistic descriptor dataset. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9180–9211.