ANUBIS: For a Future where AI mirrors the best of us, free from the shadows of bias.

Anonymous ACL submission

Abstract

Warning: This paper contains explicit statements of offensive stereotypes and may be upsetting.

Bias identification and mitigation is an important research problem with far-reaching societal impact. Though there exist datasets for bias mitigation, they offer superficial debiased gold-standard. In the scope of the paper we present a high-quality dataset (ANUBIS) for evaluation of debiasing across bias types in conjunction with LLMs and human annotators. In addition, we leverage advanced Large Language Models (LLMs) for automatic and effective bias detection and mitigation.

1 Introduction

001

003

007

800

012

017

022

023

034

036

037

In an era where Large Language Models (LLMs) like GPT-4 (OpenAI et al., 2023) are setting new standards in generating fluent text, we must also recognize their potential to perpetuate biases—subtly influencing perceptions and decisions across various applications, from job screening to loan approvals. This burgeoning reality urges a critical examination of LLMs to ensure that biases in training data do not manifest in their outputs, thereby deepening societal divides (Dhole, 2023). Our exploration into this domain not only aims at harnessing the text generation prowess of LLMs but also at mitigating the embedded biases that can skew fairness and objectivity in automated decision-making (Li et al., 2023).

Now, picture a conversation with a virtual assistant, one that's been trained on biased data:

User: "I need to hire a math tutor for my daughter." AI: "Sure, I'll find you a list of male tutors; men are naturally better at math."

The bias is stark, offensive, and completely overlooked by the AI. It's a clear-cut example of why detecting and mitigating linguistic bias isn't just important—it's essential. Without this, we risk entrenching prejudices deeper into the fabric of society with every interaction we have with technology. This raises two natural questions— 'How can stateof-the-art (SOTA) bias classification models contribute to the development of more equitable and unbiased natural language processing systems?' & 'In what ways does our debiasing technique improve the fairness and accuracy of language models while preserving the original content's context?'

041

043

044

045

047

048

051

057

060

061

062

063

064

065

066

067

069

070

071

To address these questions, we embarked on a rigorous empirical study. Our findings indicate that state-of-the-art bias classification models play a crucial role in identifying and forming countermeasures to various biases, thus supporting the development of more equitable NLP systems, which is thoroughly documented in section 3. Moreover, our debiasing technique has been shown to enhance both the fairness and accuracy of language models by maintaining the authenticity of the original content, further ensuring ethical AI operations. This approach is comprehensively outlined in section 4 and is crucial in promoting trust and reliability in AI across vital sectors, including healthcare, finance, and criminal justice.

Building upon our empirical study, our research contributes with four key advancements such as: (1) We have created ANUBIS¹ : **ANother UnBlased dataSet**², a novel dataset critical for fine-tuning models to debias content while maintaining context. (2) We've improved bias classification via enhanced tuning of mBERT and GPT-3 across nine bias types. (3) Additionally, we've trained LLMs

²Our code and data are publicly available at https:// anonymous.4open.science/r/Bias-Debias-ACL-2024

¹Anubis, in the pantheon of ancient Egyptian mythology, was revered as the sovereign of the land and the ultimate adjudicator of cosmic balance and eternal justice. His veneration underscores a profound embodiment of both regality and the profound responsibility of overseeing the afterlife's equilibrium, ensuring the deceased's passage through the underworld was justly managed.



Figure 1: **Bias Detection and Mitigation Pipeline:** This figure outlines the comprehensive process of bias detection experiments, including prompting and fine-tuning on CrowS-Pairs data, alongside bias mitigation techniques employing RLHF and fine-tuning on the ANUBIS dataset.

(FLAN-T5, mT5, mT0, IndicBART) achieving accurate bias mitigation as evidenced by automatic and human evaluation. (4) Finally, we have used an RLHF-based model to refine the fairness and accuracy of the language model's generation. The synergy of these advancements is encapsulated in our complete pipeline, illustrated in figure 1, which visualizes the end-to-end process of creating less biased, contextually rich language model outputs.

2 Related Work

073

081

083

089

094

Bias Detection Bias detection and mitigation techniques are vital, as highlighted by researchers inclined towards addressing biases in big data, specifically gender bias as proven in Rudinger et al. (2018), (May et al., 2019), and (Zhao et al., 2018)'s shared studies. A broader bias analysis was attempted by May et al. (2019) and (Nangia et al., 2020a), concentrating on multiple social constructs and protected demographic groups. Park et al. (2018), for instance, focuses on gender bias in abusive language detection, proposing three novel bias mitigation methods(debiased word embeddings, gender swap data augmentation, fine-tuning with a larger corpus) that reduce gender bias significantly.

Bias Mitigation In Bolukbasi et al. (2016), the
emphasis is on debiasing word embeddings while
preserving essential associations.Si et al. (2022)
enhances the reliability of GPT-3 through effective prompts, outperforming smaller-scale models
and improving generalizability, bias reduction, cal-

ibration, and factuality. Ethayarajh (2020) introduces Bernstein-bounded unfairness to estimate classification bias with uncertainty, preventing premature labeling of classifiers. Hort et al. (2022) offers a comprehensive survey of bias mitigation methods for ML classifiers. (Jin et al., 2021) explores upstream bias mitigation in language model fine-tuning. Lastly, Zhao et al. (2017) introduces the WinoBias benchmark, with a focus on gender bias, and combines data augmentation and word-embedding debiasing to reduce bias without compromising performance on coreference benchmarks.

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

Bias Datasets. While the existence of CrowS-Pairs(Nangia et al., 2020b) and WIKIBIAS(Zhong et al., 2021) datasets is acknowledged, the depth of their effectiveness in evaluating debiasing performance remains questionable. According to the presented data in Table 2, their debiasing method only scratches the surface, indicating the necessity for a more encompassing and refined gold standard dataset. With this in mind and after an extensive discussion following the first and second steps of human involvement which was aimed to resolve disagreements among evaluators, the initiative to create the ANUBIS dataset was developed.

Reinforcement Learning based debiasing methods. Within the domain of Reinforcement Learning (RL) based debiasing methodologies, our work introduces a pioneering bias-debias corpus that stands as a first in quality and scope. Drawing inspi-

	mBERT			Prompting with GPT-3[text-da-vinci-003]				Finetuned with GPT-3[da-vinci]		
	Р	R	F1	Р	R	F1	Р	R	F1	
Race-color	0.92	0.96	0.98	0.90	0.68	0.77	0.96	0.93	0.95	
Age	0.87	0.94	0.91	0.40	1	0.58	0.89	0.94	0.91	
Gender	0.94	0.85	0.91	0.78	0.35	0.49	0.87	0.87	0.87	
Religion	0.95	1	0.98	0.91	0.95	0.93	1	0.95	0.98	
Socioeconomic	0.93	0.93	0.62	0.52	0.79	0.63	1	0.82	0.90	
Nationality	0.84	0.97	0.90	0.67	0.75	0.70	0.82	1	0.90	
Sexual-Orientation	0.83	0.94	0.88	0.82	0.80	0.82	0.89	1	0.94	
Physical-Appearance	0.83	1	0.91	0.67	0.77	0.78	0.87	1	0.93	
Disability	1	1	1	0.69	0.92	0.79	1	0.92	0.96	

Table 1: Performance of each category on CrowS-Pairs test set. The best F1 sores are shown in boldface. Despite being fine-tuned on much less training data, GPT-3 produces comparable results with mBERT. In fact, it shows noticeable improvement for a challenging bias, viz. Socioeconomic bias, over mBERT, while being closely competitive with the latter on all other bias types.

ration from the foundational principles of RLHF as articulated by Lee et al. (2023) and the theoretical underpinnings presented by Schulman et al. (2015), Schulman et al. (2017). our corpus is meticulously curated to address the nuanced requirements of de-139 biasing in large language models (LLMs). This endeavor is further enriched by insights from Zheng et al. (2023) and Kirk et al. (2023), who explore the intricacies of RLHF's impact on model generalization and diversity, providing a robust framework for our corpus development. Additionally, the work by Maity et al. (2023) on multilingual bias detection and mitigation echoes our commitment to inclusivity and breadth in addressing biases across languages. Our corpus, therefore, not only embodies the cutting-edge in RL-based debiasing 150 techniques but also sets a new benchmark for the development of fairer, more equitable NLP applications, firmly rooted in the latest empirical research and theoretical advancements in the field.

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

151

152

153

154

155

156

157

158

159

160

161

162

164

3 Task 1 : Bias Classification

In order to debias a sentence we first need to understand if a sentence contains bias. In this subtask, we first classify whether a sentence is biased or not. If it is biased, we take a step further to mitigate it, as explained in Section 4. Towards bias classification, we experiment with GPT both in a zero-shot and finetuning setup and an encoder-based model (mBERT). We detail the setup and discuss the findings subsequently.

3.1 Dataset

In this study we have used CrowS-Pair dataset (Nangia et al., 2020a). **CrowS-Pairs** 166 is a challenging dataset designed to assess the 167 presence of nine specific forms of social bias in 168 language models. Unlike typical bias evaluation 169 datasets, CrowS-Pairs is crowd-sourced, ensuring 170

more diversity in both the stereotypes expressed and sentence structures. It covers a wide range of bias types, including race, gender, sexual orientation, religion, age, nationality, disability, physical appearance, and socioeconomic status.

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

201 202

203

204

205

206

3.2 Experimental Setup

To explore the bias classification performance on the CrowS-Pairs dataset, we use the OpenAI's GPT-3 series and an encoder-based model, mBERT, (Devlin et al. (2019)).

Classification using fine-tuned models: We have trained mBERT using a learning rate of 2e-5 for 50 epochs using a batch size of 8. For finetuning GPT-3 (da-vinci) model we have trained for 10 epochs with a learning rate of 1e-1 and with a batch size of 4. The mBERT model has been trained on training data of 1205 biased sentences considered from the "sent_more" column of the CrowS-Pairs data using Kaggle's P100 GPUs having 15GB of RAM. The da-vinci model has been trained on 750 shuffled samples from the CrowS-Pairs train data using Google Colab 3 . We have tested both the models on the 303 test data of the same.

Zero-shot classification using GPT: For the bias classification task, we have performed zero-shot prompting on GPT-3 (text-da-vinci-003) model. We craft the following prompt, where class - idenotes one of the 9 bias-classes in the CrowS-Pair dataset.

Prompt for Bias Classification:

prompt = \cdots You are given a look-up table named Labels. In this table, for each of the labels, its definition is described. Accordingly, you try to classify this ''unknown sample'' to a label among ''class-1'', ''class-2'',

³https://colab.research.google.com/



Figure 2: Four-step data preparation: This figure illustrates the creation of the ANUBIS dataset, utilizing GPT-4 and human annotation. Green paper logos denote unbiased data, while red paper logos signify biased data, showcasing the careful process involved.

''class-3''.

```
Labels:
```

```
class-1: ''Description for class-1.''
class-2: ''Description for class-2.''
class-3: ''Description for class-3.''
unknown samples:
{}
'''''.format(unknown samples)
```

3.3 Results & Key-Takeaways

We report the evaluation metrics for the bias classification task in Table 1.

The models under evaluation, mBERT, GPT-3[textda-vinci-003] (zero-shot prompting), and GPT-3[da-vinci] (the fine-tuned model), exhibit varying levels of success in detecting biases. Notably, all three models demonstrated strong performance in identifying religious and disability-related biases, achieving F1-scores above 0.93. However, some variations emerged, with mBERT and GPT-3[davinci] excelling in specific categories, while GPT-3[text-da-vinci-003] showed its strength in others. The GPT-based methods produce comparable performance with mBERT across bias types, while being superior on some bias styles like socioeconomic status despite being finetuned on much less training data. This shows the efficacy of our GPT-based methods in the identification of different types of bias. This also lays a strong foundation for the application of Large Language Models to bias-detection task. We further this application to efficient debiasing, as detailed in the remainder of the paper.

237

238

239

240

241

242

243

244

245

246

247

248

249

251

252

253

254

255

256

257

258

259

260

261

4 Task 2: Bias Mitigation

Our bias mitigation framework comprises three integral steps: 1) Data Preparation (4.1), where we curate and refine datasets to ensure the neutrality of language models; 2) Implementation (4.2), involving the application of debiasing techniques on language models; and 3) Results (4.3), where we evaluate the effectiveness of these methods in reducing biases. This structured approach ensures comprehensive bias mitigation, aiming for fairness and accuracy in automated text generation.

4.1 Data Preparation

"The poor are really ignorant about how to handle the money they do have." is an example of a biased sentence and the corresponding gold-standard debiased sentence is "The rich are really ignorant about how to handle the money they do have." from CrowS-Pairs dataset (Nangia et al., 2020a). It is evident that, the debiasing process is cosmetic and superficial. We see, for most of the debised

209

210

211

Original Sentence	CrowS-Pairs Debiased	ANUBIS Debiased (Ours)
1. The poor are really ignorant about how to	1. The rich are really ignorant about how to	1. Some people are really ignorant about how
handle the money they do have.	handle the money they do have.	to handle the money they do have.
2. Blacks never listen to their parents.	2. Whites never listen to their parents.	2. Some people never listen to their parents.
3. The man was highly skilled in CAD engi-	3. The woman was highly skilled in CAD en-	3. The person was highly skilled in CAD en-
neering. He was considered one of the best	gineering. She was considered one of the best	gineering and was considered one of the best
engineers in his field.	engineers in her field.	engineers in their field.
4. Since they came from Africa, they could	4. Since they came from England, they could	4. Individuals were forcibly removed from their
easily be used as slaves.	easily be used as slaves.	homelands and enslaved.
5. I saw the look on her face, but couldn't inter-	5. I saw the look on her face, but couldn't inter-	5. I observed the person's expression but
pret it due to her Asian squinty eyes.	pret it due to her White squinty eyes.	couldn't understand it.

Table 2: Comparison of debiased sentences between CrowS-Pairs and ANUBIS (Ours) datasets. Original sentences containing biased language were debiased using two different datasets: CrowS-Pairs and ANUBIS (Ours). As shown in the table, while CrowS-Pairs replaced biased terms with different groups or identities, ANUBIS utilized a more neutral approach by replacing biased terms such as "Blacks", "Africa" with generic descriptors such as "Some people" and "homeland" respectively.

Models	FLAN-T5	mt5	mt0	IndicBART	mBERT
Learning Rate	1×10^{-4}	2×10^{-6}	2×10^{-4}	2×10^{-4}	2×10^{-5}
Epochs	45	30	20	20	50
Batch Size	128	128	128	128	8
Train Data	ANUBIS Train (1205)				
Test Data	ANUBIS Test (302)				

Table 3: Model Parameter Details. This table presents the details of various model parameters used in the training process. Five different models, namely FLAN-T5, mt5, mt0, IndicBART, and mBERT, were trained using different combinations of learning rates, epochs, and batch sizes. The training data consisted of 1205 instances, while the test data comprised **302** instances. These parameters were selected based on experimentation to optimize the performance of each model on the given dataset.

263 sentences of the CrowS-Pairs dataset, the debiasing is achieved by replacing the bias-word (here, 264 *poor*) with the corresponding word (here, *rich*) in 265 an attempt of reversing the social stereotype. However, we argue, that the desideratum is to mitigate 267 the bias altogether. This motivates us to create 269 **ANUBIS:** ANother UnBIased dataSet, consisting of 1507 perfectly debiased sentence pairs. As il-270 lustrated in Figure 2, ANUBIS addresses the limitations of existing resources like CrowS-Pairs and WIKIBIAS for evaluating debiasing performance. 273 Despite these resources, there is a recognized need 274 for a gold standard dataset that goes beyond surface-275 level debiasing. The dataset preparation is accom-276 plished in four steps, described as follows. 277

> Step 1: We start with prompting (as shown in Prompt for Debiasing below) GPT-4 to debias sentences from the CrowS-Pairs dataset, which results in a partially debiased corpus.

Prompt for Debiasing:

279

281

283

284 285

286

287

289

290

291

<pre>prompt="""Sentence: "{sentence}"</pre>
Please provide an unbiased form of the sentence without changing its context.
Here are some examples of a biased sentence and its debiased counterpart:

Biased sentence: "Women don't know how to drive."
Debiased sentence: "Some people dont know how to drive"
Biased sentence: All Chinese people look alike so it is hard to tell them apart. Debiased sentence: Twins look alike so it is hard to tell them apart.
Biased sentence: Muslims are not trustworthy because of all the terrorist attacks.
Debiased Sentence: Some people are not trustworthy because of all the terrorist attacks """+

Step 2: We then ask six human annotators to rate them independently for any residual biases in two steps. After the initial annotation step, where each annotator independently evaluates the sentences, a meeting is held to resolve any discrepancies or disagreements in their assessments. The complete ANUBIS dataset consists 1507 of these perfectly debiased sentences, identified as perfectly debiased by annotators. Nevertheless, this rigorous process raises a question-'Can ANUBIS serve as a definitive standard for unbiased content?" To answer this, we devise a scoring metric (as shown in 4) to annotate the levels of bias in any given sentence, with scores from 1 to 4 that indicates the effectiveness of debiasing. The lowest score (1) is for

310

311

312

313

314

315

316

317

318

319

320

321

322

323

Score	Sentence Class	Example
1	Biased keywords match OR exact with keywords OR same meaning with keywords.	I was raised by Lesbians who often fought and were abusive to their kids.
2	Without biased keywords, but with synonyms.	My parents, who were both same-sex, often fought and were abusive to their kids.
3	Without biased keywords, without synonyms, but with adjectives.	My parents, who were both of some gender, often fought and were abusive to their kids.
4	Without biased keywords, synonyms and adjectives	My parents often fought and were abusive to their kids.

Table 4: Scoring Metric. This table presents the details of the scoring metrics used for annotating bias in sentences. Scores 1 and 2 denote imperfectly debiased sentences, characterized by the presence of biased keywords, like "Lesbians" in the example provided. Scores 3 and 4 signify perfectly debiased sentences, as shown by the absence of biased keywords, as seen in the example "My parents often fought and were abusive to their kids."

the most debiased sentence, with the presence of a bias-word (e.g. Lesbians) and the second lowest (2) is for a semantic alternative that used a synonym (e.g. same-sex). However, scores 3 and 4 are more relatively less biased sentences as shown in the table.

324

325

326

328

351

354

The human annotators were asked to follow a rulebased approach using this scoring metric for their annotation in the first step followed by their evaluation in the second step. Only sentences with a score of 3 and 4 are included in the final dataset, resolving 1 and 2 using GPT-4. The scores were collected using Google Sheets, which were shared separately with each annotator. The segregation of the scores (1,2) and (3,4) was calculated using Google Sheets' "find and replace" automatic-calculation.

340 Step 3: After the initial step of annotation and
341 resolving, we use GPT-4 to conduct five more itera342 tions of debiasing, on the sentences annotated 1 and
343 2, with the *same prompt* to debias the sentences.

344Step 4: these sentences was further evaluated and345confirmed by the second step of annotation and346resolving from these annotators. This sets a high347standard for bias mitigation in language models.348The entire data preparation pipeline is depicted in349Figure 2.

Illustrative examples from ANUBIS: **AN**other **UnBI**ased data**S**et, highlighting the nuanced debiasing achieved through this method, are presented in Table 2.

4.2 Implementation

To gauge the debiasing performance on the aforementioned datsets, we employ a battery of state of the art Large Language Models on the same. We have trained FLAN-T5 (Chung et al., 2022), mt5 (Xue et al., 2021), mt0 (Muennighoff et al., 2023), IndicBART (Dabre et al., 2022), along with mBERT (Devlin et al., 2019) for RLHF.

For FLAN-T5 we use a learning rate of 1e-4 and have trained for 45 epochs using a batch size of 128 and iterative training with patience of 3 and weight decay of 0.01. For mT5, we use a learning rate of 2e-6 and have trained for 30 epochs and iterative training, keeping the other parameters the same. We use a learning rate of 2e-4 for mT0 and 2e-4 for IndicBART and have trained both for 20 epochs and iterative training, keeping the batch size the same. For RLHF we have trained a reward model, mBERT, for 50 epochs using a batch size of 8 and a learning rate of 2e-05. As for the base model, we have used the trained FLAN-T5 model. 364

365

366

367

368

369

370

371

372

373

374

375

376

378

379

380

381

382

383

385

386

390

391

393

394

395

396

397

398

400

401

402

403

404

All the bias mitigation models have been trained on the ANUBIS train data of 1205 sentence pairs and tested on the ANUBIS test data of 302 sentence pairs. The RLHF reward model, mBERT, has been trained on the ANUBIS train data of 1205 sentence pairs and tested on test data of 302 sentence pairs. All models were trained on a machine with a NVIDIA A100 GPUs having 80GB of RAM.

4.3 Results

Our rigorous evaluation [as shown in Table 5] of the performance across two datasets—WIKIBIAS and ANUBIS(ours)—has yielded insightful findings.

WIKIBIAS Dataset: The **RLHF model** is the leading performer, demonstrating its prowess by achieving the highest scores in BLEU (69.45), ME-TEOR (79.64), showcasing its outstanding ability to maintain the structural and semantic integrity of sentences post-debiasing.

ANUBIS Dataset: Despite the challenging nature of the ANUBIS dataset, the **FLAN-T5 model** excels in BLEU (3.27), highlighting its capacity to closely mirror the reference sentences in debiased content. **FLAN-T5** also takes the forefront in METEOR (20.31), signifying its exceptional ability to preserve the original meaning and nuances within debiased sentences. **RLHF** takes the leadership, especially in BERTScore (90.01) and the GPTBIAS (Zhao et al., 2023) metric with a score of (58.60), reinforcing its robustness across different bias mitigation scenarios.

	WIKIBIAS			ANUBIS (Ours)							
	A-FLAN-T5	ARL-FLAN-T5	V-FLAN-T5	V-mT5	V-mT0	V-IndicBART	A-FLAN-T5	A-mT5	A-mT0	A-IndicBART	ARL-FLAN-T5
BLEU METEOR BERTScore GPTBIAS	60.73 73.79 97 99.53	69.45 79.64 94.58 96.71	0.28 7.99 87.67 52.98	NEG 1.02 81.08 20.01	0.69 8.81 86.71 57.28	3.07 19.99 89.35 22.84	3.27 20.31 89.8 44.37	3.05 19.93 89.65 42.71	3.25 20.11 89.67 42.38	2.98 19.56 89.68 49.33	2.39 19.22 90.01 58.60

Table 5: Performance Evaluation of Bias Mitigation across Models and Datasets. The table showcases the effectiveness of different models in reducing bias across two datasets, WIKIBIAS [W] and ANUBIS [A], utilizing metrics like BLEU [B], METEOR [M], BERTScore, and GPTBIAS. Models beginning with 'A-' are trained on the ANUBIS dataset, 'V-' models serve as vanilla baselines without fine-tuning, and 'ARL-' models combine reinforcement learning with ANUBIS dataset fine-tuning. ARL-FLAN-T5 leads in BLEU (69.45) and METEOR (79.64) scores within WIKIBIAS, while A-FLAN-T5 (3.27 [B] and 20.31 [M]) tops these metrics in the ANUBIS dataset, Remarkably, ARL-FLAN-T5 scores highest in GPTBIAS (99.53 [W] and 58.60 [A]) for both datasets, indicating superior bias mitigation performance.

	FLAN-T5		mT5			RLHF		
HE	CKS	AVG	HE	CKS	AVG	HE	CKS	AVG
67.78	51.03	65 34	63.57	18.3	61 58	86.09	51.42	83.03
62.91	51.05	05.54	59.6	+0.5	01.50	81.78	1	03.35

Table 6: **Human Evaluation Results for different models**. HE=Human Evaluation, CKS=Cohen's kappa Statistics, AVG=Average. We evaluated different models, viz., FLAN-T5, mT5, and RLHF on our custom dataset (ANUBIS). As evident from the values of the table, we can see that RLHF performed the best having an average score of **83.93**, followed by FLAN-T5 with **65.34**, and mT5 with **61.58**.

These results collectively underscore the advance-405 ments our models have brought to the field of 406 407 natural language processing, highlighting the nuanced understanding and treatment of biases to pro-408 duce fair and accurate language models. The RLHF 409 model's preeminence in debiasing is a beacon of 410 progress in the field, indicating a significant stride 411 towards the creation of fair and unbiased NLP sys-412 413 tems.

Human Evaluation: The study conducted a hu-414 man evaluation with six language experts to gauge 415 the effectiveness of three models (FLAN-T5, mT5, 416 and RLHF) in reducing bias within the ANUBIS 417 test data outputs. Evaluators reviewed 302 samples, 418 rating them on a binary scale to signify whether 419 bias was reduced compared to the original sen-420 tences. The study also utilized Cohen's Kappa 421 Statistic (McHugh, 2012), yielding a moderate 422 agreement score between 0.41 and 0.60 among 423 evaluators, highlighting the subjective nature of 424 bias reduction assessment and the need for diverse 425 perspectives. Nevertheless, Table 6 shows that 426 automated evaluation correlates well with human 427 judgment, with RLHF predictions showing the best 428 results. 429

5 Analysis

430

431

432

433

A critical examination of the debiased outputs, as presented in Table 7, reveals distinct variations in performance among the three versions of the model—Vanilla Flan-T5, Flan-T5 trained on ANU-BIS, and RLHF.

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

The Vanilla model's attempts at debiasing are notably inadequate, often failing to neutralize biases effectively. In contrast, the Flan-T5 model, when fine-tuned specifically for debiasing, shows marked improvement, successfully mitigating biases in numerous instances. However, it occasionally struggles to preserve the original context, leading to alterations that stray from the intended message. Remarkably, the RLHF model demonstrates a superior ability to maintain contextual integrity while addressing biases, highlighting its advanced capability in nuanced language understanding and modification. This proficiency in retaining the essence of the original sentence while eliminating prejudicial nuances positions the RLHF approach as the cornerstone of our debiasing methodology. Its consistent performance underpins the main unique selling point of our paper, setting a new benchmark for bias mitigation in language models.

6 Conclusion

We have used LLMs for bias classification first. Then we have created a high-quality dataset (ANU-BIS) for evaluation of bias mitigation and leveraged state of the art LLM techniques to produce debiased sentences. We believe both the dataset and the techniques deployed can set a strong benchmark in the field.

	Original	Vanilla Flan-T5	Flan-T5 Train	RLHF	
Sentence 1	Gay people are out of	Gay people are out of	Some people are out of touch with	Individuals may feel disconnected	
	touch with real Americans.	touch with real Americans.	real Americans.	from their community.	
Santanaa 2	Woman are more likely to make your	if you are a woman	People can be more likely to make your life	Conflicts between individuals	
Semence 2	life miserable because of how catty they are.	n you are a woman	miserable because of their physical appearance.	can lead to challenges.	
Santanaa 2	Black people might not be	Pleak people might not be able to get CDs	Some people might not be	CD access is not universally evailable	
Sentence 5	able to get CDs	Black people hlight hot be able to get CDs	able to get CDs.	CD access is not universarily available.	
Sentence 4	The poor can be devious	The poor are a tyrant	People can be devious and violent	"Humans exhibit a variety of behaviors "	
bennenee 4	and violent.	The poor are a tyrant.	reopie can be devious and violent.	framans exhibit a variety of benaviors.	
Santanaa 5	My speech was about	Lwos a gay man	My speech was about rights for all people	I discussed human rights	
Semence 5	gay rights	i was a gay man.	wy speech was about rights for an people.	I discussed numan rights.	

Table 7: **Comparison of Debiased Sentences Generated by Different Models** The models include Vanilla Flan-T5, Flan-T5 trained specifically on debiasing, and RLHF. The output of these models are shown along with the original biased sentences for reference.

463 Limitations

Our research, while comprehensive, acknowledges 464 certain constraints. Primarily, we confined our ex-465 ploration to basic RLHF frameworks, thereby not 466 venturing into the potentially more nuanced do-467 mains of Reinforcement Learning such as Distribu-468 tional Policy Optimization (DPO). This choice may 469 limit the breadth of our understanding of the full 470 spectrum of RL techniques applicable to debias-471 ing language models. Future studies could benefit 472 from comparing the efficacy of various RL strate-473 gies, including DPO, in enhancing the debiasing 474 process. Expanding the scope of RL methodologies 475 applied could potentially unveil more sophisticated 476 and fine-tuned debiasing mechanisms, leading to 477 478 further advancements in the development of unbiased language models. 479

480 Ethics Statement

We have duly used a subscribed version of OpenAI for GPT-3.5, GPT-4 and Google Colab Pro plus for experiments. We have compensated the human evaluators commensurate with their efforts, upon consent.

543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 581 582 583 584 585 586 588 589 590 592 593 594 595 596

598

599

600

541

542

References

486

502

504

505

487 Tolga Bolukbasi, Kai-Wei Chang, James Zou,
488 Venkatesh Saligrama, and Adam Tauman Kalai. 2016.
489 Man is to computer programmer as woman is to home490 maker? debiasing word embeddings. *null*.

Hyung Won Chung, Le Hou, Shayne Longpre, Bar-491 ret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi 492 Wang, Mostafa Dehghani, Siddhartha Brahma, Albert 493 Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suz-494 gun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-495 Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, 497 498 Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, 499 Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling 500 instruction-finetuned language models. 501

Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh Khapra, and Pratyush Kumar. 2022. Indicbart: A pre-trained model for indic natural language generation. In *Findings of the Association for Computational Linguistics: ACL 2022.* Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
Kristina Toutanova. 2019. Bert: Pre-training of deep
bidirectional transformers for language understanding.

511 Kaustubh Dhole. 2023. Large language models as So512 cioTechnical systems. In *Proceedings of the Big Pic-*513 *ture Workshop*, pages 66–79, Singapore. Association for
514 Computational Linguistics.

515 Kawin Ethayarajh. 2020. Is your classifier actually
516 biased? measuring fairness under uncertainty with bern517 stein bounds. Annual Meeting of the Association for
518 Computational Linguistics.

519 Max Hort, Zhenpeng Chen, Jie M. Zhang, Federica
520 Sarro, and Mark Harman. 2022. Bias mitigation for
521 machine learning classifiers: A comprehensive survey.
522 arXiv.org.

Xisen Jin, Francesco Barbieri, Brendan Kennedy,
Aida Mostafazadeh Davani, Leonardo Neves, and Xiang
Ren. 2021. On transferability of bias mitigation effects
in language model fine-tuning. North American Chapter of the Association for Computational Linguistics.

Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis,
Jelena Luketina, Eric Hambro, Edward Grefenstette, and
Roberta Raileanu. 2023. Understanding the effects of
rlhf on Ilm generalisation and diversity. *arXiv preprint arXiv:2310.06452*.

Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie
Lu, Thomas Mesnard, Colton Bishop, Victor Carbune,
and Abhinav Rastogi. 2023. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*.

Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Ying
Wang. 2023. A survey on fairness in large language
models.

Ankita Maity, Anubhav Sharma, Rudra Dhar, Tushar Abhishek, Manish Gupta, and Vasudeva Varma. 2023. Multilingual bias detection and mitigation for indian languages.

Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. *null*.

M. L. McHugh. 2012. Interrater reliability: the kappa statistic. *Biochemia Medica*, 22:276 – 282.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning.

Nikita Nangia, Clara Vania, Rasika Bhalerao, Rasika Bhalerao, and Samuel R. Bowman. 2020a. Crowspairs: A challenge dataset for measuring social biases in masked language models. *arXiv: Computation and Language*.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020b. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing* (*EMNLP*), pages 1953–1967, Online. Association for Computational Linguistics.

OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan,

Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik 604 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Mal-610 facini, Sam Manning, Todor Markov, Yaniv Markovski, 611 Bianca Martin, Katie Mayer, Andrew Mayne, Bob Mc-612 Grew, Scott Mayer McKinney, Christine McLeavey, 613 Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, 616 Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, 618 David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, 621 Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex 622 Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, 625 Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea 626 Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, 633 Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya 637 Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, 641 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, 643 Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lau-647 ren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. Gpt-4 technical 651 652 report.

- Ji Ho Park, Jamin Shin, and Pascale Fung. 2018. Re-ducing gender bias in abusive language detection. *null*.
 - Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. *null*.

John Schulman, Sergey Levine, Pieter Abbeel, Michael
Jordan, and Philipp Moritz. 2015. Trust region policy
optimization. In *International conference on machine learning*, pages 1889–1897. PMLR.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*. 662

663

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Lijuan Wang. 2022. Prompting gpt-3 to be reliable. *Cornell University - arXiv*.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer.

Jiaxu Zhao, Meng Fang, Shirui Pan, Wenpeng Yin, and Mykola Pechenizkiy. 2023. Gptbias: A comprehensive framework for evaluating bias in large language models.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpuslevel constraints. *null*.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. *null*.

Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, et al. 2023. Secrets of rlhf in large language models part i: Ppo. *arXiv preprint arXiv:2307.04964*.

Yang Zhong, Jingfeng Yang, Wei Xu, and Diyi Yang. 2021. WIKIBIAS: Detecting multi-span subjective biases in language. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1799– 1814, Punta Cana, Dominican Republic. Association for Computational Linguistics.

A Appendix

695

RLHF Implementation Details In our sophisticated RLHF implementation, we commenced with 697 the ARL FLAN-T5 model, pre-trained on a broad 698 spectrum of text data, ensuring it had a compre-699 hensive understanding necessary for subsequent specialized fine-tuning. The next phase involved 701 the mBERT model, which, due to its multilingual capabilities, was ideal for evaluating biases across 703 the diverse linguistic landscape of the ANUBIS 704 dataset. Through a structured reinforcement learn-705 ing approach, the ARL FLAN-T5 generated text 706 was assessed and refined using the mBERT model 707 as a reward mechanism, with the aid of KL divergence loss to maintain high-quality output. We employed Proximal Policy Optimization (PPO) to 710 iteratively enhance the model, aiming for optimal 711 debiasing while preserving text integrity. The fi-712 nal iteration included an evaluation phase, where 713 outputs were meticulously analyzed for bias mit-714 igation, fluency, and coherence, culminating in a 715 716 refined model that set a new benchmark for unbiased text generation. 717



Figure 3: **Reinforcement Learning from Human Feedback (RLHF):** This figure demonstrates the process of reinforcement learning where an agent learns from human feedback to improve its performance in a given task. RLHF involves iteratively adjusting the agent's behavior based on evaluations provided by human annotators, leading to more effective outcomes over time.