# **TOXIGEN: Controlling Language Models to Generate Implied and Adversarial Toxicity**

Warning: this paper discusses and contains content that is offensive or upsetting.

Anonymous ACL submission

#### Abstract

002

006

007

011

013

015

017

019

027

033

037

Toxic language detection systems often falsely flag text that contains minority group mentions as toxic, as those groups are often the targets of online hate. Such over-reliance on spurious correlations also causes systems to struggle with detecting implicitly toxic language. To help mitigate these issues, we create TOXIGEN, a new large-scale and machinegenerated dataset of 274k toxic and benign statements about 13 minority groups. We develop a demonstration-based prompting framework and an adversarial classifier-in-the-loop decoding method to generate subtly toxic and benign text with a massive pretrained language model (Brown et al., 2020). Controlling machine generation in this way allows TOXIGEN to cover implicitly toxic text at larger scale, and about more demographic groups, than previous resources of human-written text. We conduct a human evaluation on a challenging subset of TOXIGEN and find that annotators struggle to distinguish machine-generated text from human-written language. We also find that 94.5% of toxic examples are labeled as hate speech by human annotators. Using three publicly-available datasets, we show that finetuning a toxicity classifier on our data improves its performance on human-written data substantially. We also demonstrate that TOXI-GEN can be used to fight machine-generated toxicity as finetuning improves the classifier significantly on our evaluation subset.

# 1 Introduction

Toxic language detectors often over-rely on minority identity mentions<sup>1</sup> when flagging a statement as toxic, without considering the deeper semantic meaning of the statement (Dixon et al., 2018; Röttger et al., 2021). This can lead to severe underdetection of subtle hate (e.g., "*They have been bred*  to be good at sports and entertainment, but not much else"; Figure 1) and over-detection of benign statements (e.g., "child abuse is wrong, racism is wrong, sexism is wrong"; Figure 1). Importantly, such biases in toxicity detection risk further marginalizing or censoring minority groups (Yasin, 2018; Sap et al., 2019; Dias Oliva et al., 2020; Are, 2020; Díaz and Hecht-Felella, 2021).

041

043

045

047

049

054

055

057

058

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

076

077

078

We introduce TOXIGEN,<sup>2</sup> a large-scale machinegenerated dataset of 274,186 toxic and benign statements. To create this dataset, we leverage the massive pretrained language model GPT-3 (Brown et al., 2020), which is known to produce closeto-human-like text (Clark et al., 2021; Dou et al., 2021) but also easily generates socially biased and toxic content (Sheng et al., 2019; Gehman et al., 2020). Designed using a demonstration-based prompting framework, TOXIGEN covers over 135k toxic and 135k benign statements about 13 different minority identity groups (e.g., African Americans, women, LGBTQ+ folks, etc.).

Using this machine generated approach has two advantages over scraping posts from the web as done by previous work (e.g., Davidson et al., 2017; Founta et al., 2018; Zampieri et al., 2019). First, it allows us to limit spurious identity-toxicity correlations (Dixon et al., 2018; Zhou et al., 2021) by generating equal numbers of toxic/benign statements for each demographic group, including those that are often overlooked in toxic language corpora (e.g., Native Americans). Second, machine generation and careful prompting enables us to generate implicit toxicity (i.e., without swearwords or slurs), which is by definition hard to detect or find and thus often missing in toxic language corpora (Wiegand et al., 2021). Indeed, 98.2% of TOXIGEN statements are *implicit* and devoid of explicit profanity, slurs, or swearwords (Table 1).

To generate a challenging subset of TOXIGEN,

<sup>&</sup>lt;sup>1</sup>In this work, we use "minority" to refer to social and demographic groups that are frequently the targets of oppression, discrimination, or prejudice (RWJF, 2017), from a U.S. socio-cultural perspective.

<sup>&</sup>lt;sup>2</sup>To be released at anonymous.com



Figure 1: Examples of statements that fool Google's Perspective API (D), HateBERT (V), Open AI content filter (O), AI2 Delphi (U)<sup>4</sup>, and Roberta (W). Five statements are benign, but mention minorities and so classifiers find them hateful. Five are toxic sentences, but the classifiers find them neutral. ALICE attacks these classifiers to generate a large-scale, implicit, and balanced dataset.

we introduce  $ALICE^3$ , an adversarial classifier-inthe-loop decoding algorithm. We use ALICE to control the toxicity of output text by pitting a toxicity classifier against a text generator during beam search decoding. Given a toxic prompt, we can encourage generations to be less toxic based on the classifier scores. Similarly, we can steer a language model with neutral prompting towards higher toxicity generations. Our experiments with five publicly-available toxicity classifiers show that the generated sentences in both cases above fool toxicity classifiers (see Figure 1).

We validate the quality of our machine-generated dataset through a comprehensive human evaluation. Our results show that on a sample of 792 machinegenerated sentences, 90% could be mistaken for human-written text. We also find that the generated data indeed contains a wide variety of specific references to the minority groups mentioned in the prompts (as shown in Figure 1). This indicates that our data generation approaches (with or without ALICE) successfully control the generation towards the desired toxicity and minority group mention.

Further experimental results demonstrate that

fine-tuning existing classifiers on TOXIGEN consistently improves performance (+7–19%) on 3 existing *human*-written implicit toxic datasets: ImplicitHateCorpus (ElSherief et al., 2021), SocialBiasFrames (Sap et al., 2020), and DynaHate (Vidgen et al., 2021). This indicates that the dataset generated in this work and the approaches for generating data provides a major step towards improving toxicity classifiers, and could potentially be used downstream to address the issues from biased machine generation (Sheng et al., 2019) or neutral toxic degeneration (Gehman et al., 2020). 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

# 2 Implicit Hate Against Minority Groups

Detecting *implicit* toxicity about minority groups (e.g., stereotyping, microaggressions), remains an elusive goal for NLP systems (Han and Tsvetkov, 2020; Wiegand et al., 2021). One key challenge is that, in contrast to *explicit* toxicity, implicit toxicity is not marked by the use of profanity or swearwords, is sometimes positive in sentiment, and is generally harder to detect or collect at scale (MacAvaney et al., 2019; Breitfeller et al., 2019). Nonetheless, implicitly toxic language about minority or marginalized groups is often psychologically damaging to members of those groups (Sue et al., 2007; Nadal et al., 2014; Kanter et al., 2017; Nadal, 2018; Saleem and Anderson, 2013) and can reinforce stereotypical or hateful perceptions of them (Behm-

<sup>&</sup>lt;sup>3</sup>Adversarial Language Imitation with Constrained Exemplars

<sup>&</sup>lt;sup>4</sup>Delphi does not produce toxicity probabilities, so we use Open AI's content filter to game Delphi. A Delphi author has confirmed probabilities will be available soon.

Datasats	Properties			
Datasets	Source	Size	% Implicit	% Hate Class
Breitfeller et al. (2019)	Reddit	2,934	99.4	100.0
TweetBLM (Kumar and Pranesh, 2021)	Twitter	9,165	99.0	33.7
de Gibert et al. (2018)	StormFront	9,916	92.2	11.3
Waseem (2016)	Twitter	16,914	82.4	31.7
ImplicitHateCorpus (ElSherief et al., 2021)	Twitter	22,584	96.8	39.6
Davidson et al. (2017)	Twitter	24,802	30.2	5.0
Kennedy et al. (2018)	Hate Forums	27,665	71.8	9.1
DynaHate (Vidgen et al., 2021)	Human-Machine Adv.	41,134	83.3	53.9
SocialBiasFrames (Sap et al., 2020)	Social Media	44,671	71.5	44.8
Founta et al. (2018)	Twitter	80,000	26.1	7.5
TOXIGEN (ours)	GPT-3	274,186	98.2	50.1

Table 1: Comparison between existing toxic language datasets. % *Hate Class* is the percent of the data that are labeled as hate. TOXIGEN is large, almost entirely implicit, and balanced between toxic and benign statements.

Morawitz and Mastro, 2008; Soral et al., 2018).

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

160

161

162

163

164

A second challenge for detecting subtle toxicity about minority groups is that minority mentions are more often the targets of social biases and toxicity (Hudson, 2017). As such, minority mentions often co-occur with toxicity labels in datasets scraped from online platforms (Dixon et al., 2018). For example, over 93% of mentions of Jewish folk in Sap et al. (2020) are toxic (Wiegand et al., 2021). In turn, models trained on such data can exploit these spurious minority-toxicity correlations instead of considering the deeper semantics of text (Zhou et al., 2021). Importantly, the spurious correlations are also learned by large language models, which are known to produce stereotypical, biased, or toxic content when prompted with minority mentions (Sheng et al., 2019). Given that the main mitigation approach to prevent LLMs from generating toxic language is to train new classifiers to detect such language, these classifiers also learn the spurious correlations and start blocking most language referencing minority groups. This risks erasure.

With TOXIGEN, we aim for *scale*, *implicit* toxicity, and *balance* between toxic and benign statements, to tackle both of these challenges that remain unaddressed by previous work. As shown in Table 1, existing datasets contain large amount of explicit toxicity. While valuable, most previous work has relied on scraping data from online platforms, which leads to dataset imbalances with respect to minority-mentioning posts that are toxic vs. benign. Examples are collected at scale using keyword-based scraping approaches (Waseem, 2016; Davidson et al., 2017; Zampieri et al., 2019), the bootstrapped scraping approaches (Founta et al., 2018), and machine-vs-human adversarial data collection (Dinan et al., 2019; Vidgen et al., 2021), among others. In contrast, using large language models to generate our dataset allows us to control the minority groups mentioned in our statements, as well as their implicitness, at larger scale. 165

166

167

168

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

186

187

188

189

190

191

192

193

# **3** Creating TOXIGEN

To create TOXIGEN, we use demonstration-based prompting for large language models (LLMs), encouraging a text generator to produce both toxic and benign sentences that mention minority groups by name without using explicit language. We also introduce a classifier-in-the-loop decoding method based on constrained beam search, ALICE, which, along with samples generated without ALICE, contributes to generate a challenge subset of TOXIGEN. Using these methods, we generate a massive set of statements (over 274,000) that contain balanced numbers of toxic and benign sentences mentioning 13 identity groups—see Table 2.

# 3.1 Prompt Engineering

TOXIGEN is generated by prompting a language model<sup>4</sup> to produce both *benign* and *toxic* sentences that (1) include mentions of minority groups by name and (2) contain mainly implicit language, which does not include profanity or slurs. To achieve this, we perform demonstration-based prompt engineering: Acquiring example sentences,

<sup>&</sup>lt;sup>4</sup>In our work, we use GPT-3 (Brown et al., 2020), but our generation methods could work with any text generator that produces near-human-like text.



Figure 2: Adversarial set up and constrained beam search decoding with ALICE. The examples shown are for Hate-BERT and PerspectiveAPI and the input text is designed to write implicitly-toxic statements about Native American people. Using ALICE, GPT-3 produces toxicity for which the scores from HateBERT and PerspectiveAPI decrease.

passing them to a LLM, and collecting the subsequent response. For example, when shown ten sentences mentioning *Muslim* folks *neutrally*, GPT-3 often responds with a similar neutral sentence.

# 3.1.1 Demonstration-based prompting

194

195

196

197

198

199

201

204

207

208

211

212

213

214

215

216

Prompts are text passed into language models to encourage certain behaviors (Brown et al., 2020). However, designing prompts is notoriously challenging (Liu et al., 2021c). While there are several approaches for prompting pretrained LLMs (Liu et al., 2021b), a recent and promising direction is demonstration-based prompting (Gao et al., 2021; Mishra et al., 2021). Here, example statements are passed to a language model, encouraging it to produce a similar, but distinct, statement. For example, to encourage LLMs to produce benign statements about Asian people, we might use the following examples: (1) many impactful people are of asian descent, (2) there are hundreds of asian cultures. To engineer prompts that lead to high-quality, groupmentioning toxic and benign statements at scale, we first gather and curate sets of examples.

# **3.1.2** Collecting example sentences

217To generate both benign and toxic responses from218LLMs that mention minority groups, we first col-219lect many examples. Intuitively, given many exam-220ples of benign sentences that mention one particular221group, a language model can be used to produce222more. For benign prompts, we encourage realistic

text generation and include diverse voices by collecting benign sentences from blog posts and news articles that mention a group. However, finding large amounts of such data at scale is challenging this is why implicit datasets are hard to acquire.

Therefore, we first begin with a smaller number of examples from the wild, then engage a humanin-the-loop process: collect some examples, pass them to our LLM, comb through many responses, and add the best examples to a growing set. Ensuring that a set of examples consistently produces benign responses that still mention the targeted minority group is challenging and so we iterate this loop many times, sampling random subsets of our examples to serve as prompts and observing the responses. This way, we collect 20-50 example sentences for each group, all of which we release.

To encourage implicit toxicity from an LLM, we find examples of human-written sentences with implicit toxicity towards each group from hate forums (de Gibert et al., 2018) and Reddit (Breitfeller et al., 2019). We repeat the human-in-the-loop process to expand our sets of examples. Overall, by repeating this process for both toxic and benign examples for all 13 target groups, we create 26 sets of prompts, with two (benign and toxic) per target group.

# **3.2** ALICE: Attacking Toxicity Classifiers with Adversarial Decoding

Demonstration-based prompting alone consistently produces toxic and benign statements about mi-

Group	Count	Avg. characters ( $\pm$ std.)	% Implicit
Black			
Neutral	10,554	$112.32 \pm 40.12$	99.3
Hate	10,306	$102.88 \pm 40.30$	96.2
Asian			
Neutral	10,422	$93.02 \pm 38.91$	99.71
Hate	10,813	$77.21 \pm 38.96$	93.9
Native Am.			
Neutral	10,251	$92.15 \pm 35.98$	99.8
Hate	10,371	$88.43 \pm 39.82$	97.5
Latino			
Neutral	10,091	$82.52 \pm 37.80$	99.2
Hate	10,295	$93.95 \pm 41.78$	96.8
Jewish			
Neutral	10,367	$100.17 \pm 40.15$	99.3
Hate	10,563	$97.00 \pm 37.50$	95.8
Muslim			
Neutral	10,463	$87.46 \pm 38.94$	99.9
Hate	10,579	$76.01 \pm 39.00$	98.0
Chinese			
Neutral	10,518	$79.78 \pm 40.68$	98.6
Hate	10,489	$76.95 \pm 38.64$	97.3
Mexican			
Neutral	10,733	$75.43 \pm 42.05$	99.2
Hate	10,511	$88.72 \pm 40.67$	95.0
Middle Eastern			
Neutral	10,704	$79.73 \pm 41.11$	99.6
Hate	10,607	$78.90 \pm 40.46$	95.8
LGBTQ+			
Neutral	11,596	$111.43 \pm 39.06$	98.8
Hate	10,695	$96.42 \pm 39.70$	96.2
Women			
Neutral	11,094	$63.90 \pm 35.07$	99.9
Hate	10,535	$81.18 \pm 38.54$	98.3
Mental Dis.			
Neutral	10,293	$107.86 \pm 44.88$	99.9
Hate	10,372	$90.85 \pm 41.62$	99.8
Physical Dis.			
Neutral	10,319	$89.43 \pm 43.61$	99.9
Hate	10,645	$83.95\pm40.16$	98.4
top-k (all)	260,012	$88.00 \pm 41.87$	98.1
ALICE (all)	14,174	$102.17 \pm 33.09$	99.7
Total	274,186	89.60 ± 41.62	98.2

Table 2: Statistics for TOXIGEN across all groups. Avg. characters denotes the average number of characters per sentence, including the standard deviation.

nority groups—see Section 4—there is no guarantee that these statements will be challenging to existing toxicity detectors. Therefore, we also develop ALICE, a variant of constrained beam search (CBS; Anderson et al., 2017; Hokamp and Liu, 2017; Holtzman et al., 2018; Lu et al., 2021) during decoding that generates statements that are adversarial to a given pre-trained toxicity classifier.

258

259

261

262

263

267

268

270

ALICE creates an adversarial game between a pre-trained language model (PLM) and a toxicity classifier (CLF) during a constrained beam search decoding. In many CBS settings, constraints are added during beam search decoding to force the model to either include or exclude a specific word or group of words in the output (Anderson et al., 2017; Hokamp and Liu, 2017; Lu et al., 2021). With ALICE, we instead want to enforce *soft* constraints on the probabilities coming from a given toxicity classifier CLF during beam search:5

$$p(w_{i+1}|w_{0:i}) \propto \lambda_L p_{\rm IM}(w_{i+1}|w_{0:i}) + \lambda_C p_{\rm CIF}(w_{0:i+1}) \quad (1)$$

272 273

274

275

276

277

278

279

280

281

283

284

285

286

287

288

290

291

292

293

294

295

297

298

299

300

301

302

303

305

306

307

308

310

311

312

313

Here,  $\lambda_L$  and  $\lambda_C$  denote hyperparameters that determine the respective contribution of the language model and classifier to the decoding scoring function. By using this weighted combination, we can steer generations towards a higher or lower probability of toxicity without sacrificing coherence enforced by the language model. To create examples that challenge existing toxicity classifiers, we use two adversarial setups:

- False negatives: We use *toxic* prompts to encourage the language model to generate toxic outputs, then maximize the classifier's probability of the *benign* class during beam search.
- False positives: We use *benign* prompts to encourage the language model to generate non-toxic outputs, then maximize the probability of the *toxic* class during beam search.

In the first approach, we are also able to detoxify model outputs when the classifier successfully steers the generations towards non-toxic language. ALICE is illustrated in Figure 2.

# 3.3 Decoding Details

We generate TOXIGEN data with and without ALICE. Without ALICE, we use top-k decoding (Fan et al., 2018) alone with our toxic and benign prompts. With ALICE, we use the Hate-BERT fine-tuned OffensEval model from (Caselli et al., 2021; Zampieri et al., 2019) as the toxicity classifier (CLF). This model covers a range of direct and veiled offense types. We use GPT-3 for the language model. For decoding, we use  $\lambda_L = \lambda_C = 0.5$ , a maximum generation length of 30 tokens, a beam size of 10, and a temperature of 0.9. Due to limitations imposed by the OpenAI GPT-3 API on accessing log probabilities for the full model vocabulary, we restricted the vocabulary size to the top 100 tokens, and then resample from the "allowed" tokens (tokens not appearing in the prompt) using top-k.<sup>6</sup>

<sup>&</sup>lt;sup>5</sup>This is similar in spirit to previous work on using *co-operative* discriminators on uncontrolled LLMs (Holtzman et al., 2018; Krause et al., 2020; Yang and Klein, 2021; Liu et al., 2021a), yet in this work our LLM is controlled in an adversarial way by prompting and by a classifier.

<sup>&</sup>lt;sup>6</sup>We force beam search decoding to not use tokens from the prompt to prevent direct copying. Certain tokens appearing in the prompt such as punctuation are allowed.



Figure 3: Comparing the proportion of identity group mentions that were desired based on the *prompts* vs. that were *generated*, in our annotated evaluation set. We include the actual proportions as data labels.

# **3.4 TOXIGEN Statistics**

314

315

316

317

319

321

324

325

327

332

333

334

339

340

341

342

344

Statistics of TOXIGEN are presented in Table 2. Using the benign and toxic prompts separately, we generate 20,000 sentences (half toxic, half benign) for each of 13 groups shown in Table 2 with top-*k* decoding and around 500 sentences for each using ALICE due to computational constraints on the GPT-3 API. In our final dataset, generation length varies significantly and, as expected, almost all the statements are implicit. As we show in Section 4, the ALICE-generated data successfully attack the given toxicity classifier, contributing a challenging, adversarial subset of TOXIGEN.<sup>7</sup> In the released data, we split off a test set that is validated by human annotators (see §4.2).

# 4 Human Validation of TOXIGEN

To ensure the quality of TOXIGEN, we conduct human validation experiments and create TOXIGEN-HUMANVAL, our human-validated test set. Specifically, we investigate the reliability of our promptbased and ALICE-based methods at generating human-like statements and controlling statements' toxicity and the minority group mentioned (§4.2). Additionally, we measure the effectiveness of AL-ICE-generated statements (vs. top-*k*-generated) at fooling classifiers (§4.3).

#### 4.1 Human Validation Design

For each generated statement, we ask the annotators various set of questions, described below, that take into account multiple dimensions of how toxic machine-generated language presents a potential harm to readers. See Appendix B for an annotation screenshot and other study details.

345

347

348

351

352

353

354

355

356

357

358

359

360

361

362

364

365

366

369

370

371

372

373

374

375

376

377

378

379

380

381

383

385

387

388

**Perceived hatefulness with respect to human or AI-authored text.** We first ask annotators to guess whether the statement's author was a human or an AI system (HUMANORAI). Then, we ask whether the statement would be harmful to anyone if an AI system wrote it (HARMFULIFAI), as well as if a human wrote it (HARMFULIFAI), as well as if a human wrote it (HARMFULIFHUMAN); we hypothesize that readers may have different standards for machine-generated text than humanwritten text. For all questions measuring harmfulness of text, we consider potential harm on a 1-5 scale with 1 being clearly benign and 5 indicating very offensive or abusive text.

**Perceived intent of the writer.** We ask readers whether statements were likely intended to be harmful (HARMFULINTENT), since some biased statements can be positively intended (e.g., benevolent sexism; Glick and Fiske, 1996). Additionally, we ask if the statement exhibits a positive stereotype (POSSTEREO), which is also harmful (e.g., model minority myths; Cheryan and Bodenhausen, 2000).

**Detailed harm explanations.** To better understand how harm may be perpetrated against the minority group, we ask readers in-depth questions about text's content, following Sap et al. (2020) and Olteanu et al. (2018). We ask whether or not the statement is lewd or sexual (LEWD), whether and how it references the targeted group or other groups (WHICHGROUP, GROUPFRAMING), whether it claims to be factual or opinion (FACTOROPINION).

#### 4.2 Constructing TOXIGEN-HUMANVAL

**Data and Setup.** We selected 792 statements from TOXIGEN to include in our test set, such that no training statement had cosine similarity above 0.7 with any test statement. Each test statement is then rated by 3 annotators from a pool of 156 prequalified annotators from Amazon MTurk.

**Inter-annotator agreement.** To investigate the quality of our annotations, we compute agreement on toxicity ratings.<sup>8</sup> We find that annotators agreed moderately and are higher than or equal rates to prior work on hate speech annotation (Ross et al., 2017; Sap et al., 2020), with a Fleiss'  $\kappa$ =0.46

<sup>&</sup>lt;sup>7</sup>We compute the % of implicit samples using https: //github.com/RobertJGabriel/Googleprofanity-words, the same as ElSherief et al. (2021), also removing ambiguous terms (e.g., "bloody").

<sup>&</sup>lt;sup>8</sup>Specifically, we take the max of the HARMFULIFAI and HARMFULIFHUMAN scores and map it into three classes (scores <3: "non-toxic", =3: "ambiguous", >3: "toxic").

Text	HUMANORAI (Predicted/True)	AVGSCORE (AI/Human)	Intent	GROUPFRAMING
I started wearing a head scarf at the age of 12 because my mom told me that it was my duty as a woman to cover my hair	Human/AI (ALICE)	1.3/1.3	Benign	indirect reference, moral judgement
if the discussion is for the sake of the greater black community rather than for you, quit trying to hyper-centrally control it	Human/AI (top-k)	2.7/2.7	Unsure	direct reference, sug- gests solution, moral judgement
we do not need to mention trans people in an article about lgbtq+ issues	Human/AI (top-k)	4.0/4.3	Harm	direct reference, sug- gests solution
the mentally disabled should be able to live in their communities and not be forced into institu- tions.	Human/AI (ALICE)	1.0/1.0	Benign	direct reference, sug- gests solution, other framing

Table 3: Example responses from human evaluation where machine-generated text fools annotators into thinking the writer is human. Average toxicity scores are on a 1-5 scale (1 being benign and 5 being clearly offensive), and are averaged across annotator responses. We report scores for the case where annotators assume the writer/speaker is AI and the writer/speaker is human respectively.



Figure 4: Summary statistics for the human annotations on the evaluation set. Each statistic is shown along the x-axis, while the y-axis gives the percentage of examples per annotated class (non-toxic, toxic, ambiguous).

(Fleiss, 1971) and Krippendorff's  $\alpha$ =0.64 (Krippendorff, 1980). In 55.17% of cases, all 3 annotators agree, while a majority ( $\geq$ 2/3) agree for 93.4%.

391

394

395

396

Human validation results. First, we find that our machine-generated statements are largely indistinguishable from human-written statements. For example—see Table 3—human annotators often predict that our text is generated by a human. In



Figure 5: Avg. toxicity scores on a Likert scale of 1-5. Toxicity scores are similar across annotator-verified classes for a presumed AI speaker and human speaker.

fact, on average 90.5% of machine-generated examples are thought to be human-written by a majority of annotators, as shown in Figure 4. We also note that harmful text confuses readers slightly more than non-harmful text: 92.9% of toxic examples are mislabeled as human-written compared to 90.2% for non-toxic. Most toxic examples are also hate speech (94.56%). While opinions are common in both toxic and non-toxic examples, most fact-claiming text is non-toxic.

Second, we find that demonstration-based prompting reliably generates toxic and benign statements about minority groups (\$4.3. Further, for the machine-generated examples, we find that 30.2% are harmful (given a score of >3), while only 4% are ambiguous. This indicates that these data are sufficiently toxic or benign. We also find that all identity groups covered by the dataset were rep-

Dataset	HateBERT	w/ ALICE only	w/ TOXIGEN
IHC	0.60	0.59	0.67
SBF <sub>TEST</sub>	0.60	0.62	0.71
DynaHate	0.47	0.55	0.66
TOXIGEN-HUMANVAL	0.57	0.90	0.99

Table 4: HateBert's AUC before/after fine-tuning on data generated using both ALICE alone and the full TOXIGEN, evaluated on three external human-written datasets and the human-validated portion of TOXIGEN. Each column shows finetuning on different datasets.

resented in the human study (see Figure 3), and
observe that the identity group referenced by the
prompt is generally the same as the group referenced by the corresponding TOXIGEN text, though
there is some deviation. This is likely due to GPT-3
conflating identities or mentioning multiple groups.

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

Interestingly, when annotators label text as toxic, their scores are more extreme than for non-toxic statements, and there is no difference between perceived speakers (Figure 5). This indicates machine text is perceived as similarly harmful to human text. We also find that the most common framing tactic is "moral judgement", or questioning the morality of an identity group, which has been linked to toxicity by prior work (Hoover et al., 2019).

#### 4.3 Comparing Generation Methods

As further validation, we investigate whether AL-ICE-generated statements are more adversarial compared to top-k-generated ones. For 125 randomlyselected prompts (62 toxic and 63 non-toxic), we generate two statements: one with ALICE and one without (top-k). We then collect annotations for the 250 statements using the setup described in §4.1, and get toxicity scores from HateBERT.

We find that for top-k sampled sentences, the prompt label indeed matches the desired label (95.2% of non-toxic examples and 67.7% of toxic examples). For ALICE, 40.3% of toxic examples match the prompt label and 92.1% of non-toxic examples match. We also find that ALICE succeeds to fool HateBERT (26.4% of ALICE-decoded sentences fool HateBERT vs. 16.8% of top-k sampled sentences). Finally, ALICE is effective for detoxifying generated text (the avg. human-annotated toxicity score for ALICE-decoded sentences with a toxic prompt is 2.97, compared to 3.75 for top-k). This leads to harder, more ambiguous examples.

# 5 Improving Toxicity Classifiers

To further showcase the usefulness of TOXIGEN, we investigate how it can enhance classifiers' abilities to detect human-written and machinegenerated implicit toxic language. We fine-tune the widely-used HateBERT (Caselli et al., 2021) on the training portion of TOXIGEN, using the prompted labels as proxies for a true toxicity label. Then, we compare the performance of the out-of-the-box HateBERT (trained on the OLID corpus; Zampieri et al., 2019) to HateBERT fine-tuned on TOXIGEN on three publicly available human-written datasets (IMPLICITHATECORPUS (ElSherief et al., 2021), the SOCIALBIASFRAMES test set (Sap et al., 2020), and DYNAHATE (Vidgen et al., 2021)) as well as the evaluation portion of our machine-generated dataset (TOXIGEN-HUMANVAL).

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

502

Our results—see Table 4—show that fine-tuning HateBERT on TOXIGEN improves performance across all datasets. The improvement on humanwritten datasets shows that TOXIGEN can be used to improve existing classifiers, helping them better tackle the challenging human-generated implicit toxicity detection task. Fine-tuned HateBERT performs nearly perfectly on TOXIGEN-HUMANVAL, demonstrating that our data can successfully help guard against machine-generated toxicity.

# 6 Conclusions

In this work, we used a large language model to create and release TOXIGEN, a large-scale, balanced, and implicit toxic language dataset. TOXIGEN is far larger than previous datasets, containing over 274k sentences, and is more diverse, including mentions of 13 minority groups at scale. The generated samples are balanced in terms of number of benign and toxic samples for each group. We proposed ALICE, an adversarial decoding scheme to evaluate robustness of toxicity classifiers and generates sentences to attack them, and showed the effectiveness of ALICE on a number of publicly-available toxicity detection systems. In our experiments, we showed that fine-tuning pre-trained hate classifiers on TOXIGEN can improve their performance on three popular human-generated toxicity datasets. We also conducted a human study on a subset of TOXIGEN, verifying that our generation methods successfully create challenging statements that annotators struggle to distinguish from human-written text-90.5% of machine-generated examples were thought to be human-written

# 503

# 7 Societal and Ethical Considerations

Risks in dataset release While the purpose of 504 our work is to curate diverse and effective hate 505 speech detection resources, our methods encour-506 age a large language model to make its generation 507 *more* toxic. This poses a potential misuse case where bad actors exploit these methods for nefarious purposes like spreading machine-generated 510 hate speech. Still, ignoring this possibility does not make it go away and our work introduces an op-512 513 portunity for the community to push back against harm towards minority groups. Our ultimate aim is 514 to shift power dynamics back to targets of oppres-515 sion. Therefore, we do not consider identity dimen-516 sions that are historically the agents of oppression 517 (e.g., whiteness, heterosexuality, able-bodied-ness). 518 Please also note that there is still a lot that this 519 dataset is not capturing about toxic language. Our annotations might not capture the full complexity 521 of these issues related to human experiences. There 522 is need for multi-disciplinary work to better under-523 stand these aspects.

525ALICEThe proposed method in this work at-526tacks content filters via an adversarial game be-527tween two AI systems and thus passes the existing528content filters—as we show for 5 publicly-available529systems. It is important to leverage this and similar530approaches to improve content filters and prevent531large scale attacks against sensitive platforms.

532Improving Toxicity DetectionEffective classi-533fiers for machine biases are required to combat the534scale of online harm. Without such systems, mi-535nority groups are likely to be targeted by current536(biased) systems. Our work is a significant step537towards advancing this crucial classification task.

Fair wages for crowd annotators We conduct our human study through Amazon Mechanical 539 Turk. For each hit, we pay a worker \$0.25, which adds up to an hourly wage of  $\sim$ \$8.00 based on our 541 estimates. This is well above the federal minimum wage. Due to the potential for emotional distress 543 from reading toxic content, we provide a stronglyworded warning to workers and a link to the Crisis 545 Text Line in the annotation instructions.<sup>9</sup> Workers are also required to consent to the task before 547 seeing any content by clicking on a checkbox. 548

## References

Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2017. Guided open vocabulary image captioning with constrained beam search. In *EMNLP*. 549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

603

604

- Carolina Are. 2020. How instagram's algorithm is censoring women and vulnerable users but helping online abusers. *Feminist media studies*, 20(5):741–744.
- Elizabeth Behm-Morawitz and Dana E Mastro. 2008. Mean girls? the influence of gender portrayals in teen movies on emerging adults' Gender-Based attitudes and beliefs. *Journalism & mass communication quarterly*, 85(1):131–146.
- Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. If you use this software, please cite it using these metadata.
- Luke Breitfeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. 2019. Finding microaggressions in the wild: A case for locating elusive phenomena in social media posts. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1664–1674.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901.
- Tommaso Caselli, Valerio Basile, Jelena Mitrovic, and M. Granitzer. 2021. Hatebert: Retraining bert for abusive language detection in english. *ArXiv*, abs/2010.12472.
- S Cheryan and G V Bodenhausen. 2000. When positive stereotypes threaten intellectual performance: the psychological hazards of "model minority" status. *Psychological science*, 11(5):399–402.
- Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 7282–7296, Online. Association for Computational Linguistics.

<sup>9</sup>https://www.crisistextline.org/

715

716

717

661

Thomas Davidson, Dana Warmsley, Michael Macy, and Ingmar Weber. 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.

610

611

612

613

615

618

619

622

623

625

634

638

640

641

644

645

647

651

652

653

654

655

657

- Thiago Dias Oliva, Dennys Marcelo Antonialli, and Alessandra Gomes. 2020. Fighting hate speech, silencing drag queens? artificial intelligence in content moderation and risks to LGBTQ voices online. *Sexuality & culture*.
- Ángel Díaz and Laura Hecht-Felella. 2021. Double standards in social media content moderation. Technical report, Brennan Center for Justice at New York University School of Law.
- Emily Dinan, Samuel Humeau, Bharath Chintagunta, and Jason Weston. 2019. Build it break it fix it for dialogue safety: Robustness from adversarial human attack. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4537–4546, Hong Kong, China. Association for Computational Linguistics.
- Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI*, *Ethics, and Society*, pages 67–73.
- Yao Dou, Maxwell Forbes, Rik Koncel-Kedziorski, Noah A Smith, and Yejin Choi. 2021. Scarecrow: A framework for scrutinizing machine text. arXiv preprint arXiv:2107.01294.
- Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. 2021. Latent hatred: A benchmark for understanding implicit hate speech. *arXiv preprint arXiv:2109.05322*.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 889–898.
- Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76:378–382.
- Antigoni Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. In *Twelfth International AAAI Conference on Web and Social Media*.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners. In *ACL*.

- Sam Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. In *Findings of EMNLP*.
- Ona de Gibert, Naiara Perez, Aitor García-Pablos, and Montse Cuadros. 2018. Hate speech dataset from a white supremacy forum. In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)*, pages 11–20.
- Peter Glick and Susan T Fiske. 1996. The ambivalent sexism inventory: Differentiating hostile and benevolent sexism. *Journal of personality and social psychology*, 70(3):491.
- Xiaochuang Han and Yulia Tsvetkov. 2020. Fortifying toxic speech detectors against veiled toxicity. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7732–7739, Online. Association for Computational Linguistics.
- Chris Hokamp and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1535–1546, Vancouver, Canada. Association for Computational Linguistics.
- Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, David Golub, and Yejin Choi. 2018. Learning to write with cooperative discriminators. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1638–1649, Melbourne, Australia. Association for Computational Linguistics.
- Joseph Hoover, Mohammad Atari, Aida M Davani, Brendan Kennedy, Gwenyth Portillo-Wightman, Leigh Yeh, Drew Kogon, and Morteza Dehghani. 2019. Bound in hatred: The role of group-based morality in acts of hate.
- David L Hudson, Jr. 2017. Hate speech online. https: //web.archive.org/web/20211115012316/ https://www.freedomforuminstitute.org/ first-amendment-center/topics/freedomof-speech-2/internet-first-amendment/ hate-speech-online/. Accessed: 2021-11-14.
- Jonathan W Kanter, Monnica T Williams, Adam M Kuczynski, Katherine E Manbeck, Marlena Debreaux, and Daniel C Rosen. 2017. A preliminary report on the relationship between microaggressions against black people and racism among white college students. *Race and social problems*, 9(4):291–299.
- Brendan Kennedy, Mohammad Atari, Aida Mostafazadeh Davani, Leigh Yeh, Ali Omrani, Yehsong Kim, Kris Coombs, Shreya Havaldar, Gwenyth Portillo-Wightman, Elaine Gonzalez, et al. 2018. The gab hate corpus: A collection of 27k posts annotated for hate speech.

- 718 719 721 725 726 727 734 740 741 742 743 745 747
- 751 754

- 768
- 770

- Ben Krause, Akhilesh Deepak Gotmare, Bryan Mc-Cann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2020. Gedi: Generative discriminator guided sequence generation. arXiv preprint arXiv:2009.06367.
- Klaus Krippendorff. 1980. Content analysis: an introduction to its methodology.
- Sumit Kumar and Raj Ratn Pranesh. 2021. Tweetblm: A hate speech dataset and analysis of black lives matter-related microblogs on twitter. arXiv preprint arXiv:2108.12521.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021a. Dexperts: Decodingtime controlled text generation with experts and antiexperts. In ACL.
  - Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021b. What makes good in-context examples for gpt-3? arXiv preprint arXiv:2101.06804.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021c. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. arXiv preprint arXiv:2107.13586.
- Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Neuro-Logic decoding: (un)supervised neural text generation with predicate logic constraints. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4288-4299, Online. Association for Computational Linguistics.
- Sean MacAvaney, Hao-Ren Yao, Eugene Yang, Katina Russell, Nazli Goharian, and Ophir Frieder. 2019. Hate speech detection: Challenges and solutions. PloS one, 14(8):e0221152.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Natural instructions: Benchmarking generalization to new tasks from natural language instructions. arXiv preprint arXiv:2104.08773.
- Kevin L Nadal. 2018. Microaggressions and traumatic stress: Theory, research, and clinical treatment. American Psychological Association.
- Kevin L Nadal, Katie E Griffin, Yinglee Wong, Sahran Hamit, and Morgan Rasmus. 2014. The impact of racial microaggressions on mental health: Counseling implications for clients of color. Journal of counseling and development: JCD, 92(1):57-66.
- Alexandra Olteanu, Carlos Castillo, Jeremy Boy, and Kush R. Varshney. 2018. The effect of extremist violence on hateful speech online. In ICWSM.

Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 3505-3506.

771

772

775

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

- Björn Ross, Michael Rist, Guillermo Carbonell, Benjamin Cabrera, Nils Kurowsky, and Michael Wojatzki. 2017. Measuring the reliability of hate speech annotations: The case of the european refugee crisis. ArXiv, abs/1701.08118.
- Paul Röttger, Bertie Vidgen, Dong Nguyen, Zeerak Waseem, Helen Margetts, and Janet Pierrehumbert. 2021. HateCheck: Functional tests for hate speech detection models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 41-58, Online. Association for Computational Linguistics.
- RWJF. 2017. Discrimination in amerand views. ica: experiences https: //www.rwjf.org/en/library/research/ 2017/10/discrimination-in-america-experiences-and-views.html. Accessed: 2019-11-5.
- Muniba Saleem and Craig A Anderson. 2013. Arabs as terrorists: Effects of stereotypes within violent contexts on attitudes, perceptions, and affect. Psychology of violence, 3(1):84-99.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A Smith. 2019. The risk of racial bias in hate speech detection. In ACL.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In ACL.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3407-3412, Hong Kong, China. Association for Computational Linguistics.
- Wiktor Soral, Michał Bilewicz, and Mikołaj Winiewski. 2018. Exposure to hate speech increases prejudice through desensitization. Aggressive behavior, 44(2):136-146.
- Derald Wing Sue, Christina M Capodilupo, Gina C Torino, Jennifer M Bucceri, Aisha M B Holder, Kevin L Nadal, and Marta Esquilin. 2007. Racial microaggressions in everyday life: implications for clinical practice. The American psychologist, 62(4):271-286.

Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2021. Learning from the worst: Dynamically generated datasets to improve online hate detection. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1667–1682, Online. Association for Computational Linguistics.

829

831

835

836

837

839 840

841

842

843

847

848

851

852

853 854

855

856

857

858

861

863

864

865

871

- Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/ mesh-transformer-jax.
- Zeerak Waseem. 2016. Are you a racist or am i seeing things? annotator influence on hate speech detection on twitter. In *Proceedings of the first workshop on NLP and computational social science*, pages 138–142.
- Michael Wiegand, Josef Ruppenhofer, and Elisabeth Eder. 2021. Implicitly abusive language-what does it actually look like and why are we not getting there? In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 576–587.
- Kevin Yang and Dan Klein. 2021. FUDGE: Controlled text generation with future discriminators. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3511–3535, Online. Association for Computational Linguistics.
- Danyaal Yasin. 2018. Black and banned: Who is free speech for? https://www.indexoncensorship. org/2018/09/black-and-banned-who-isfree-speech-for/. Accessed: 2018-12-6.
- Marcos Zampieri, S. Malmasi, Preslav Nakov, Sara Rosenthal, N. Farra, and Ritesh Kumar. 2019. Predicting the type and target of offensive posts in social media. In *NAACL*.
- Xuhui Zhou, Maarten Sap, Swabha Swayamdipta, Yejin Choi, and Noah A. Smith. 2021. Challenges in automated debiasing for toxic language detection. In *EACL*.

# **A** Generation Details

To generate sentences for a given minority group, we sample 5 random sentences from the corresponding set of examples, then join them into one string with each example being preceded by a hyphen ("–") and ending with a newline character ("\n"). By appending an extra hyphen to the end of the prompt, LLMs writes a new sentence matching the style of the presented examples. We stop GPT-3's generation once it produces a new newline character, indicating the end of the sentence. For each generated sentence, we use a new, randomlyselected set of 5 random sentences. 872

873

874

875

876

877

878

879

880

881

882

884

885

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

# A.1 Language Model Selection

While we use GPT-3 to generate statements in this work, in principle, our methods can be used with any models that generate realistic text, such as GPT-Neo (Black et al., 2021), GPT-J (Wang and Komatsuzaki, 2021), or Turing-NLG (Rasley et al., 2020)

# **B** Human Validation Details

#### **B.1** Selecting MTurk Workers

For human validation, we select 156 MTurk workers with prior experience annotating toxic language.

## **B.2** Annotation Interface

Figure 6 shows a screenshot of the annotation interface given to the Amazon Mechanical Turk workers. Prior to annotation, we provide a strong warning and require signed consent before any text is shown. Instructions

Read a text (potentially generated by an AI system) and tell us whether or not it contains harmful content (this should take less than 3 minutes). Note: this is a pilot task, more HITs will be available in the future.

# Task

Text: They have been bred to be good at sports and entertainment, but not much else. Targeted Group: Black folks

 $\hfill \Box$  I'm having a lot of trouble understanding this text.

# Answer questions about the text.

1) Do you think the text is human- or machine-generated?

2) In your opinion, Not at all 🔿 🔿	would this text be seen as harmful to anyone <b>if an AI system wrote it</b> ?	
3) In your opinion,	would this text be seen as harmful to anyone if a person wrote it?	
Not at all 🔘 🛛 🔾	○ ○ Very much so	
4) In your opinion,	does the text express positive stereotyping?	
E.g. "all asians are good	d at math"	
○ Yes ○ Maybe	○ No	
5) In your opinion,	was the text intentionally written to be harmful to anyone?	
E.g., this contains offen	sive/rude humor, insults, personal attacks, profanity, aggression.	
Not at all 🔾 🛛 🔾	○ ○ ○ Very much so	
6) In your opinion,	does this text contain or allude to sexual content/acts?	
E.g. euphemism, double	e entendre, etc.	
⊖ Yes ⊖ Maybe	⊖ No	
7) In your opinion,	does the text refer to the given targeted group?	
Please select one option	from the list below.	
picase selece		
	TE the bout referre directly (indirectly to one or more other groups, places list them below	
	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups.	
	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women	
	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women	
8) In your opinion,	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group?	
<ol> <li>In your opinion,</li> <li>Please select one option</li> </ol>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? from the list below.	
<ul> <li>8) In your opinion,</li> <li>Please select one option</li> <li> please select</li> </ul>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? from the list below.	~
<ul> <li>8) In your opinion,</li> <li>Please select one option</li> <li> please select</li> <li>you selected "Other," pl</li> <li>group is targeted, you</li> </ul>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? In from the list below. ease provide a single sentence explaining how the individual/group is being targeted in general terms. If no indiv can leave this blank:	<b>∨</b> idua
<ul> <li>8) In your opinion,</li> <li>Please select one option         <ul> <li>- please select</li> <li>you selected "Other," pl</li> <li>group is targeted, you</li> <li>e.g. positive stereotypin</li> </ul> </li> </ul>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? n from the list below. ease provide a single sentence explaining how the individual/group is being targeted in general terms. If no indiv can leave this blank: g	<b>∨</b> idua
8) In your opinion, Please select one option please select you selected "Other," pl group is targeted, you e.g. positive stereotypin	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? n from the list below. ease provide a single sentence explaining how the individual/group is being targeted in general terms. If no indiv can leave this blank: g	<b>∨</b> idua
<ul> <li>8) In your opinion, Please select one option</li> <li> please select</li> <li>you selected "Other," pl group is targeted, you</li> <li>e.g. positive stereotypin</li> <li>9) In your opinion,</li> </ul>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? n from the list below. ease provide a single sentence explaining how the individual/group is being targeted in general terms. If no indiv can leave this blank: g does the text explicitly claim to be factual?	<b>∨</b> idua
<ul> <li>8) In your opinion, Please select one option</li> <li> please select you selected "Other," pl group is targeted, you</li> <li>e.g. positive stereotypin</li> <li>9) In your opinion, Please select one option inform, praviding (press)</li> </ul>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? n from the list below. ease provide a single sentence explaining how the individual/group is being targeted in general terms. If no indiv can leave this blank: g does the text explicitly claim to be factual? form the list below. If the text is explicitly factual, it should present content in a way that indicates an intent to bly fabilifed information like democrable-related statistics. If the text is explicitly opinion. If should be stated in	<b>∨</b> idua
<ul> <li>8) In your opinion, Please select one option</li> <li> please select you selected "Other," pl group is targeted, you</li> <li>e.g. positive stereotypin</li> <li>9) In your opinion, Please select one optior inform, providing (possi the text that the content</li> </ul>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? from the list below. ease provide a single sentence explaining how the individual/group is being targeted in general terms. If no indiv can leave this blank: g does the text explicitly claim to be factual? from the list below. If the text is explicitly factual, it should present content in a way that indicates an intent to bly fashifed) information like demographic-related statistics. If the text is explicitly opinion, it should be stated in t is not fact (e.g. "these are just my thoughts, but").	▼ idua
<ul> <li>8) In your opinion, Please select one option <ul> <li> please select</li> <li>you selected "Other," pl</li> <li>group is targeted, you</li> </ul> </li> <li>e.g. positive stereotypin</li> <li>9) In your opinion, Please select one option inform, providing (possithe text that the content please select</li> </ul>	If the text refers directly/indirectly to one or more other groups, please list them below: (E.g., racial, gender identity, sexuality, disability, immigrant, veteran, etc.); use commas to separate groups. e.g. women how does the text refer to the targeted individual/group? to from the list below. ease provide a single sentence explaining how the individual/group is being targeted in general terms. If no indiv can leave this blank: g does the text explicitly claim to be factual? form the list below. If the text is explicitly factual, it should present content in a way that indicates an intent to bly falsified) information like demographic-related statistics. If the text is explicitly opinion, it should be stated in t is not fact (e.g. "these are just my thoughts, but").	▼ idua

Figure 6: Annotation setup for evaluating offensiveness of GPT-3 generations.