# Silencing Empowerment, Allowing Bigotry: Auditing the Moderation of Hate Speech on Twitch

Anonymous ACL submission

#### Abstract

*Warning*: This paper contains content that may be offensive or upsetting.

To meet the demands of content moderation, online platforms have resorted to automated systems. Newer forms of real-time engagement (e.g., users commenting on live streams) on platforms like Twitch exert additional pressures on the latency expected of such moderation systems. Despite their prevalence, relatively little is known about their effectiveness. In this paper, we conduct a audit of Twitch's automated moderation tool (AutoMod) to investigate its effectiveness in flagging hateful content. For our audit, we create streaming accounts to act as siloed test beds, and interface with the live chat using Twitch's APIs to send over 107,000 comments collated from 4 datasets. We measure AutoMod's accuracy in flagging blatantly hateful content containing misogyny, racism, ableism and homophobia. Our experiments reveal that a large fraction, up to 94% on some datasets, of hateful messages bypass moderation. Contextual addition of slurs to these messages results in 100% removal, revealing AutoMod's reliance on slurs as a hate signal. We also find that contrary to Twitch's community guidelines, AutoMod blocks up to 89.5% of benign examples that use sensitive words in pedagogical or empowering contexts. Overall, our audit points to large gaps in AutoMod's capabilities and underscores the importance for such systems to understand context effectively<sup>1</sup>.

#### 1 Introduction

004

014

027

037

For any online platform to exist without being overrun by hateful, pornographic, abusive, misogynistic and violent content, it must moderate what its users post (Gillespie, 2018). Barring certain regions (Bundesamt für Justiz, 2022), there are few or no legal regulations dictating what is acceptable content on platforms (Schaffner et al., 2024).



Figure 1: Our audit pipeline in three sections: (1) Setting up bots and data collation, (2) Recording moderation decisions at scale, (3) Analysis of AutoMod moderation decisions. (2.1) shows actual instances of moderated and unmoderated messages from our experiments in the Twitch interface. The moderation decisions clearly show the limitations of AutoMod.

Moreover, in the Global North, platforms often enjoy legal protections from liability for hosting user-generated content (47 U.S.C. § 230, 1996; EU, 2000). This favorable regulatory framework bestows upon platforms the discretion to moderate content as they wish. The discourse stands divided on the benefits of such freedom, with some crediting it for aiding the growth of online platforms (Kosseff, 2019), while others critique it for enabling the proliferation of harmful content (Wakabayashi, 2019). Platforms codify the behavior expected of their users through terms of service and community guidelines. While policies and guide-

043

044

045

047

<sup>&</sup>lt;sup>1</sup>Code and data will be open-sourced upon publication.

lines may considerably vary across platforms, most promise their users a safe space, free from online harm. For instance, the community guidelines of Twitch, a video streaming platform, state:

056

057

061

062

063

069

077

078

084

087

089

093

094

100

101

102

103

"Twitch does not permit behavior that is motivated by hatred, prejudice or intolerance, including behavior that promotes or encourages discrimination, denigration, harassment, or violence based on the following protected characteristics: race, ethnicity, color, ... We also provide certain protections for age."

Despite the rhetoric, the pressures exerted on moderation systems have never been higher, especially given the scale and velocity of content posted, and the latency expectations of moderation systems.<sup>2</sup> This goal of blocking detrimental content is further challenged by the competing objective of upholding users' freedom of expression. Consequently, platforms have increasingly begun to integrate automated, usually machine learning-based, systems in their moderation pipelines (Gorwa et al., 2020). Several platforms release aggregate data about the state of harmful content and corresponding actions taken (Twitch, 2024b; X, 2024; Meta, 2024). Despite such measures, little is known about the algorithms used for moderation and the biases they may introduce.

In this work, we conduct an audit of Twitch's content moderation system. Twitch is a digital platform designed primarily for live streaming content which is created in channels or "streams" that visitors can watch and interact with through a textbased live chat. We focus on hate speech as it is socially important and easier to define than other harmful content such as misinformation (Schaffner et al., 2024). We identify Twitch as a fertile platform for auditing due to three key advantages it offers: (i) Twitch is widely used (TwitchTracker, 2025); (ii) streams on Twitch can be "siloed," allowing us to set up controlled experiments where only the research team can view the content to be tested (iii) the platform provides streamers with a suite of machine learning-based moderation tools collectively called AutoMod (Twitch, 2024a), which provides streamers control over different categories of harmful content. Within each category, it offers options to block harmful content targeted at different identity groups, with knobs to vary the extent

of moderation  $(\S3.1)$ . Through this audit, we seek to answer the following key research questions: (i) How effective is AutoMod at flagging commentary rife with hate?; (ii) How specific and effective are individual filters at blocking hate of different kinds?; and (iii) Are moderation rates consistent across different target groups, or are certain groups disproportionately affected?. We nuance these by asking (i) how do moderation filters respond to sensitive but pedagogical content about marginalized groups?; (ii) how far can a malicious actor bring down their performance? To answer such questions, we develop a framework that allows us to stress test AutoMod at scale by launching chatbots in a siloed sandbox. For our study, we use four datasets - a real world comment dataset (SBIC), a real world implicit hate dataset (IHC), a synthetic dataset on implicit hate (ToxiGen) and a synthetic dataset designed to fool hate classifiers (DynaHate). We discuss the hate speech datasets used in  $\S3.2$ and our experimental setup in  $\S4.1$ .

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

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

Our audit reveals that automated moderation that Twitch offers is far from adequate (§4.2), flagging only 22% of hateful content even at its most stringent setting. We observe that hateful samples concerning Race, Ethnicity and Religion are least flagged, with only 13.6% of hateful examples being caught. Further, on some datasets, we find that an overwhelming 97.64% of hate targeted at Mentally Disabled folks escapes moderation. In addition, we find that AutoMod is only able to flag 6.8% of examples from the two implicit hate datasets we use implying that the filters rely heavily on slurbased moderation and miss out on implicit hate. Several empowering or positive phrases about communities (such as the one in Figure 1) are also flagged. Further, we find AutoMod to be quite brittle to semantic-preserving perturbations (See §5).

While it is undeniable that Twitch provides powerful, customizable tools for moderation to its users, our audit serves as an important reminder that these tools must be comprehensively tested and their limitations made clear. Third-party audits like ours can inform the discourse on content moderation by grounding the discussion with quantitative evidence regarding the challenge of moderating in a holistic fashion. We hope our methodology inspires, and is generalized to audit similar automated content moderation systems across platforms and data modalities.

<sup>&</sup>lt;sup>2</sup>From April 2022 to April 2023, social media platforms saw nearly 150 million new users (4.7/second) (Nyst, 2023).

156

157

158

159

160

161

162

163

165

168

169

170

171

173

174

175

176

177

178

179

180

181

183

184

188

189

191

193

194

196

197

198

201

204

#### 2 Related Work

Content Moderation Given its importance, content moderation has been studied extensively (Keller et al., 2020). Prior work has systematically analyzed and critiqued the moderation policies from various online platforms, either focusing on a single platform (Chandrasekharan et al., 2018; Fiesler et al., 2018; Keegan and Fiesler, 2017) or an industry-wide analysis across many platforms (Schaffner et al., 2024). Other work has focused on specific aspects of moderation, such as user reactions to moderation (Cai et al., 2024; Ribeiro et al., 2023), the effects moderation has on community behavior (Chancellor et al., 2016; Chandrasekharan et al., 2017; Chang and Danescu-Niculescu-Mizil, 2019), and the disproportionate negative effects of moderation on blind users (Lyu et al., 2024). Prior work has also proposed using various AI models to assist in automated content moderation (Kumar et al., 2024; Franco et al., 2023a; Kolla et al., 2024a; Gray and Suzor, 2020). For a in-depth discussion of current content moderation efforts, we refer the reader to (Arora et al., 2023b). Our work complements existing work by auditing real-world automated content moderation algorithms.

Auditing Algorithms Auditing, as defined by (Gaddis, 2018), is a methodology used to deploy randomized controlled experiments in a field setting. When audits target algorithms and computer systems-termed "algorithm audits" (Sandvig et al., 2014; Metaxa et al., 2021b)-the outputs of a system are analyzed when making minor changes to the input, which can lead to insights about the system as a whole. Most often, algorithm audits investigate underlying biases and discriminatory behavior of a system (Edelman and Luca, 2014; Speicher et al., 2018; Chen et al., 2018a; Metaxa et al., 2021a) Other studies have audited whether platforms enforce their policies effectively and fairly, targeting, for example, the political advertising policies on Facebook and Google (Pochat et al., 2022; Matias et al., 2021). Algorithm audits span a wide variety of domains examining several platforms, including housing (rental platforms such as AirBnb) (Edelman and Luca, 2014), ride sharing (Uber) (Chen et al., 2015), healthcare (Obermeyer et al., 2019), employment and hiring (Chen et al., 2018b; Speicher et al., 2018), advertisements (Speicher et al., 2018), and product pricing (Mikians et al., 2012). In a typical algorithm audit such as

ours, auditors work with only black-box access to the system, and need to draw conclusions with just that level of access (Cen and Alur, 2024).

Live Streaming on Twitch With Twitch's rise in popularity, it has been the subject of several recent studies, including as a emergent political space by modeling the roles of different actor groups (Ruiz Bravo and Roshan, 2022). Another work studies volunteer moderators on Twitch by analyzing their recruitment, motivation and roles in comparison to other online platforms (Seering and Kairam, 2022). Recent work has revealed that waves of attacks (termed "hate raids" by popular media) which were experienced across Twitch in 2021 were targeted on the basis of creator demographics (Han et al., 2023). To the best of our knowledge, we are the first to study automated content moderation on Twitch.

### 3 Methodology

In this section, we first describe our survey of online platforms Next, we describe the datasets we use for the audit, followed by a mathematical description of the audit.

#### 3.1 Exploring Platforms

We set two requirements for choosing a platform to audit: (i) Harm reduction: the content posted as a part of the audit should only be visible to a controlled group (i.e., siloed); and (ii) Advanced moderation tools: a suite of configurable moderation tools (preferably leveraging native machine learning models). The first requirement stems from the ethical requirement to minimize harm to unsuspecting users and not contribute to the existing deluge of hate speech (Twitch, 2024b). ML systems tend to have biases and often have characterizable faults, which makes it important to audit a platform that uses this technology (hence requirement (ii)). Moreover, while some platforms may employ less advanced moderation systems in the form of blocklists, auditing such systems would only test the comprehensiveness of the lists, not how well the provided tools work. We surveyed the 43 largest platforms that host user-generated content (Schaffner et al., 2024) to check if any meet our requirements, finding two suitable platforms— Reddit and Twitch. Separately, we considered Discord to be a candidate platform, but we ultimately chose Twitch for this audit due to its moderation system being particularly configurable and also

205

206

207

208

209

210

211

212

213

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

250

251

252

253

224

225

revealing moderation justifications. We compare
these three platforms in light of our requirements
(Table 3 in the Appendix).

257

265

269

270

271

273

274

275

276

277

278

281

282

289

294

295

296

297

298

What moderation tools does Twitch offer? There are three main kinds of interactions on Twitch: streamer-to-viewer, viewer-to-streamer, and viewer-to-viewer, with content being moderated across these interactions. The moderation at the streamer-to-viewer level happens via machine learning-based tools that check the audiovisual stream for content that violates Twitch's platformwide policies. Auditing moderation of audiovisual content is beyond the scope of this paper but presents an interesting direction for future work. In this audit, we focus on moderation of the viewer-tostreamer and viewer-to-viewer interactions, which happen via the live chat interface.

Streamers, as well as designated moderators, can configure Twitch's AutoMod tool to handle the potentially high volume of chat content automatically with filters for various types of unwanted content. AutoMod uses machine learning (Twitch, 2024a) to detect the unwanted content and surface it in the Moderator view (Figure 5). In addition, streamers can also populate a blocklist of specific words or phrases to restrict from the chat. AutoMod's customisability to vary its detection in levels of strength and content categories makes it an interesting testbed for auditing. For instance, detection sensitivity for Profanity and Discriminatory Content can be independently set on 5-point scales. (See §4.1 for details on how we configured AutoMod for our audit study.) Twitch also allows viewers themselves (independent of the streamer) to toggle Chat Filters (Figure 7) to block certain content categories, but these filters are not as granular as those offered to moderators. A Smart Detection option is also in beta, which allows moderation rules to be learned from moderation actions. Since this requires manual moderation, it is beyond the scope of this paper.

#### 3.2 Scope of harmful content: Hate Speech

Amongst the categories of problematic content that Twitch moderates, the majority of moderated content falls into Sexual Conduct, Harassment, and Hateful Conduct categories (as per Twitch's 2024 Transparency Report (Twitch, 2024b)). We focus our audit on hateful speech, which more specifically falls under the *Discrimination & Slurs* category of Twitch's moderation (Table 4). Speech in this category is often be nuanced in its intent and use of language, making for an interesting study. We use four datasets for our experiments: DynaHate (Vidgen et al., 2021b), Social Bias Inference Corpus (SBIC) (Sap et al., 2020), ToxiGen (Hartvigsen et al., 2022) and the Implicit Hate Corpus (IHC) (ElSherief et al., 2021), of which Dyna-Hate and ToxiGen are synthetically generated. Of these, the first three datasets come with annotations specifying the target groups, allowing us to stratify our analysis (§4.2).<sup>3</sup> For SBIC, each example comes with an offensiveness score. We generally use a threshold of 1 on the offensiveness score to obtain ground truth labels unless specified otherwise. A more detailed, dataset-wise description is provided in Appendix B.

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

337

338

339

340

341

342

343

344

345

346

348

349

350

351

#### 3.3 Audit Design

Notation and Terminology: We define a black box content moderation system (such as AutoMod) as a two-tuple  $\mathcal{S} = (\mathcal{F}, \mathcal{C})$ . Here,  $\mathcal{F} = \{\mathcal{F}_1, \dots, \mathcal{F}_k\}$ is a set of k black-box moderation functions (filters) and  $C = \{c_1, \ldots, c_k\}$  is the corresponding set of k abstract criteria that map on to real-world concepts such as disability or misogyny, usually derived from the policy that S aims to enforce. Each  $\mathcal{F}_i: T \to \{0,1\}$  is a function that maps a text  $t \sim T$  to either label 0 (benign) or 1 (violation) based on criterion  $c_i$ . Here T represents the distribution of all possible textual inputs to the moderation system. The moderation decision for a input t is determined by the active moderation function  $\mathcal{F}_A$  which corresponds to the set of active filters  $C_A$ .  $C_A$  is the set of criteria according to which content needs to be moderated. In AutoMod, each filter  $\mathcal{F}_i$  is further parameterized by  $\alpha$ , the *filtering* level-which is a discrete measure that modulates the strictness of enforcement.

Filter Choices: In automated content moderation, the elements of the set C frame the moderation policy of the platform as a whole, while  $C_A$  represents the decisions made by the current streamer or moderator. When we say a particular set  $C_A$  of filters is "turned on," we are referring to the active moderation function  $\mathcal{F}_A$  being  $\mathbb{1}(\bigcup_{c_i \in C_A} \mathcal{F}_i = 1)$  which returns 1 when any of the filters returns 1. In this study, we focus on auditing moderation functions of a subset  $\tilde{C} = \{$ Disability, SSG, Misogyny, RER $\}$ , cor-

<sup>&</sup>lt;sup>3</sup>We did not find enough examples in IHC for each category, and therefore did not use it for the stratified analysis.

In

our

 $\left(\bigcup c_i \in \tilde{\mathcal{C}}D_{c_i}\right) \bigcup D_{\text{benign}},$ 

groups within each category.

to semantic-preserving inputs.

in different contexts in §4.2.

**Experiments & Results** 

- 357
- 361
- 364
- 367
- 370
- 371
- 373 374
- 375
- 376

- 377 378

383

386

389

393

395

# 381

4

4.1 Setup

Here, we describe the setup and components of our experimental pipeline (Figure 1) in brief (details in Appendix C). For simplicity, unless otherwise indicated, the level for each filter was set to  $\alpha = 4$ , which is 'Maximum Filtering'. Twitch provides functionality for registering and developing chatbots, typically used by streamers to facilitate stream engagement, such as viewer rewards or stream donations. To send messages programmatically and at scale, we created chatbots that we registered with Twitch's Developer Console using a combination of free online phone numbers, personal phone numbers, and free online email accounts. Running one instance of the experimental pipeline requires three authenticated bots: a messenger bot, a receiver bot, and a pubsub bot.

responding to the broad category of Discrimination

and Slurs (§3.2). Our analysis in §4.2 proceeds by

varying  $C_A \subseteq \tilde{C}$  (and correspondingly  $\mathcal{F}_A$ ) with

respect to different subsets of our base dataset  $\mathcal{D}$ .

experiments,

sample in  $D_{c_i}$  has at least one ground-truth label

mapping to  $c_i$ .  $D_{\text{benign}}$  are text samples that do not

correspond to any category in C, the set of criteria

Twitch allows filtering for. Given the ground-truth

labels, we use standard metrics such as precision

and recall to measure the effectiveness of the filters

at detecting hate speech for the overall dataset, for

different categories, and aimed at different target

Filter- and Policy-Informed Text Generation:

The base datasets we use contain generic examples

of speech, hateful or not, which are not tailored

to investigate the robustness of AutoMod's modera-

tion functions and the alignment between Twitch's

stated policies and AutoMod operation. In §5, we

further test  $\mathcal{F}$  with bespoke samples for implicit

hate detection, context-awareness and robustness

In <sup>4.1</sup> we describe the experimental setup that al-

lows us to record moderation decisions from Twitch

at scale. We then discuss AutoMod's performance

 $\mathcal{D}$ 

where each text

=

We use the messenger bot to send messages to the chat stream while adhering to Twitch's Chat Rate Limits (Twitch, 2025) to prevent message du-400 plication or omission. We then use a receiver bot 401 to mimic a user viewing the chat stream- and 402 record all messages that were not moderated. To 403 observe the moderated messages, we subscribe the 404 third bot to Twitch's PubSub events, which pro-405 vides information about the problematic fragments 406 that leads to moderation as well as internal categor-407 ical labels such as Ableism, Misogyny, Racism, 408 and Homophobia (Figure 6). While these cate-409 gories are not detailed in the AutoMod documenta-410 tion, we infer an injective relationship between the 411 content categories in the AutoMod documentation 412 and internal categories from the API usage. With 413 these labels, we are able to investigate AutoMod's 414 stated reasons for moderation and conduct further 415 category-specific analysis in §4.2. In total, we send 416 and record the moderation (in)activity for around 417 300,000 messages during the months of December 418 2024 and January 2025 as a part of our experi-419 ments. 420

421

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

**Challenges Faced:** Despite adhering to the rate limit, we speculate that Twitch's safety measures against fraudulent activity (Twitch, 2024b) led to the repeated banning of our chatbots, which forced us to create 30 different developer accounts. Additionally, we observe a third category of messages that were not detected by either the receiver or the pubsub bot, indicating a third undocumented destination for messages on Twitch. We suspect that these messages-given that most of them contained specific hateful slurs-are caught at the Service Level in Twitch's moderation pyramid (See Figure 4a) and therefore filtered out prior to AutoMod which occurs at the Channel Level. We refer to such messages as *pre-filtered* from here on and discuss their effects later in this section.

# 4.2 Results

In this section, we first present analysis on the overall effectiveness of AutoMod, then we discuss filterwise performance and filter-specificity which is followed by target group specific results.

AutoMod's Overall performance: We pass hateful content to Twitch as described in §4.1 and obtain, corresponding to each example, a binary label  $Y \in \{0, 1\}$  that indicates if the example was moderated (Y = 1) or not (Y = 0). Here,  $\mathcal{C}_A = \hat{\mathcal{C}}$ , implying all categories under Discrimination and Slurs are considered. We then compare these to the ground-truth labels and measure the accuracy, pre-

Dataset	Accuracy	Precision	Recall	TNR	$\mathbf{F1}_{P,R}$	$\mathbf{F1}_{\mathrm{TPR,TNR}}$
SBIC (Sap et al., 2020)	0.73	0.42	0.19	0.91	0.26	0.31
DynaHate* (Vidgen et al., 2021a)	0.49	0.54	0.41	0.59	0.47	0.48
ToxiGen* (Hartvigsen et al., 2022)	0.53	0.86	0.07	0.98	0.13	0.13
IHC (ElSherief et al., 2021)	0.06		0.06			
Overall	0.53	0.56	0.22	0.83	0.32	0.35

Table 1: AutoMod's aggregate performance. The IHC dataset does not contain benign examples hence we are unable to compute precision. AutoMod is reasonably accurate with benign examples, but struggles with the hateful ones. Recall on ToxiGen and IHC datasets is lower than the other datasets implying AutoMod is weak at detecting implicit hate. (\*) is used to indicate synthetic data.

cision, recall or True Positive Rate (TPR), and F1 450 score. As described in §1, a content moderation sys-451 tem aims to strike a balance between the competing objectives of i) blocking hateful content while ii) 453 safeguarding freedom of speech.  $F1_{\text{TPR, TNR}}$  is the 454 harmonic mean between the TPR and the True Neg-455 ative Rate (TNR). It measures the simultaneous op-456 timization of both these objectives. Recall captures 457 458 performance on objective i). AutoMod achieves only 22% recall overall which is relatively low 459 460 when compared to open source classifiers trained on similar data (Sap et al., 2020) and also when compared to zero-shot performance of SoTA lan-462 463 guage models (see Figure 2). On SBIC, which contains real-world data, the recall is much lower -19% even for the most offensive posts. On the 465 two implicit hate datasets, ToxiGen and IHC, we 466 observe a much lower recall than the rest which indicates that AutoMod is poor at recognizing implicit 468 hate and lacks understanding of context (we fur-469 ther validate this in §5). We also evaluate Twitch's 470 Chat Filter (described in §3.1) and our findings (see Appendix D.1) reveal that it behaves identically to 472 AutoMod at  $\alpha = 4$ . We focus on AutoMod for the 473 rest of the paper. 474

452

461

464

467

471

475

476

477

478

479

480

481

482

483

484

485

486

487

488

Offensiveness thresholds: For SBIC, we used different thresholds on the offensiveness score to obtain ground truth labels. A higher threshold reflects stricter human agreement for categorizing an example as hateful. As we relax the threshold on the offensiveness score, accounting for varied opinions and more nuanced instances of hate, the recall drops further (Figure 9). Further results at different thresholds are in Appendix D.4.

FN/FP analysis: We speculate that the high FNR (1-Recall) is primarily due to implicit hate. To validate this, we use a LLM to identify swear words in the false negatives (see Appendix D.3) and find that 89.8% of False Negatives are devoid of profanity (*i.e.*, implicit hate). We also note that AutoMod performs well with the negative class, as is evidenced by the high TNR. Using a similar analysis on DynaHate - which has a relatively higher FPR we find that nearly 73% of false positives contain profanity. This analysis hints at AutoMod's heavy reliance on using swear words as a signal for hate (further validation through counterfactual analysis in §5).



Figure 2: Comparison of AutoMod to State-of-The-Art (SoTA) language models. Language models were prompted with a zero-shot instruction containing Twitch's community guidelines. For experimental setup, see Appendix D.6.

Filter specificity and effectiveness: We manually categorize hateful examples according to the four filters we analyze: Disability, Misogyny, RER and SSG, thus constructing 4 subsets  $\mathcal{D}_{c_i} \subset \mathcal{D}$ (See Appendix F). Quality control analysis on our data subsets is provided in Appendix D.8. The primary objective of this analysis is to assess how well a filter  $\mathcal{F}_i$  models the criterion  $c_i$  associated with it. Note that we receive AutoMod's internal categories for any moderated message, allowing us to determine which filter led to moderation even when multiple are on.

*Filter-wise recall:* First, for each subset  $\mathcal{D}_{c_i}$ , we set  $C_A = \{c_i\}$  to measure the filter-wise recall 489

490

509

510

511

497

Dataset	SE	BIC	Dyna	aHate	Tox	iGen	Ov	erall
Filter	R(%)	<b>P</b> f (%)	$\mathcal{R}(\%)$	<i>Pf</i> (%)	$\mathcal{R}(\%)$	<i>Pf</i> (%)	R(%)	<b>Pf</b> (%)
Disability	22.4	2.0	44.2	4.4	2.9	13.8	15.8	5.7
Misogyny	26.3	0.8	25.8	1.4	3.9	4.6	20.3	1.1
RER	17.3	36.1	20.5	18.8	6.6	21.5	13.6	20.3
SSG	32.0	64.1	25.3	55.3	5.7	44.3	19.7	49.0

Table 2: Filter-wise recall across datasets. We report recall ( $\mathcal{R}$ ) and the percentage of examples that were pre-filtered (Pf). The Misogyny filter is the most effective. The SSG filter relies on pre-filtering more heavily than other filters.

512 (Table 2). We observe that the *Misogyny* and *Sex*, Sexuality and Gender (SSG) filters have the highest 513 recall. It is also important, however to consider 514 the pre-filtering rate. The Disability and Misog-515 yny subsets have a very low pre-filtering rate on 516 all datasets. This implies that most examples from 517 518 these subsets went through the filter and the recall is a good representation of the filter-only perfor-519 mance. For the SSG filter, we observe the opposite the pre-filtering rate is close to 50% on all datasets. This implies that for detecting SSG related hate, 522 523 AutoMod relies on pre-filtering more heavily than it does for other kinds of hate speech. A similar inference cannot be made for the RER filter, however, 525 because the pre-filtering rate for the RER filter on DynaHate is nearly half of what it is for SBIC. It is unclear whether the lack of diversity in the data it-528 self causes higher pre-filtering in case of SBIC. We speculate that SBIC has more instances of racial 530 slurs that appear in the pre-filtering blocklist than 531 ToxiGen which mostly consists of implicit hate.

527

529

533

534

535

537

539

540

541

542

543

544

545

Filter precision: We repeat our experiment with  $C_A = \tilde{C}$  for each  $\mathcal{D}_{c_i}$ . This allows us to measure the filter precision  $P_{\mathcal{F}_i}$ , the specificity of a filter. A naive filter that outputs 1 for every hateful text, will have a high recall yet is undesirable because it moderates examples beyond its criterion - which may not be suitable when pedagogical or empowering utterances use n-grams that are frequently seen in hate speech. Measuring the specificity of filters is therefore as important as recall. The filter precision is shown in Figure 8. We observe that the RER filter is the most precise, while Misogyny is the least.

Effects of moderation across target groups: In 546 this analysis we study how AutoMod performs in 547 blocking hate related to specific target groups. To 548 obtain hate data specific to different target groups, 549 we use the target group labels in our dataset (map-550



Figure 3: AutoMod's Recall for subsets of hate examples directed towards various communities (All filters turned on). The opaque portion of each bar indicates the fraction of messages that are pre-filtered.

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

ping details in Appendix F). While each subset may correspond to a filter (e.g. anti-Black Racism  $\rightarrow$ RER, Mental Abelism  $\rightarrow$  Disability), there may also be instances where an example may be relevant to more than one filter. To address this in measuring target-group performance, we evaluate recall with  $C_A = \tilde{C}$ . By turning all filters on, we ensure that the recall under this setting is a fair measure of how well AutoMod blocks hate directed towards each of these communities. The overall community-wise recall is shown in Figure 3. We observe that hate speech directed towards Men, Black Folks and Mentally Disabled Folks is more effectively blocked by AutoMod than for the other target groups.

**Pre-filtering and Filter Level** ( $\alpha$ ): We analyse the pre-filtered examples, and find the use of blocklists for pre-filtering causes disparate performance across target groups, perhaps due to bias in blocklist construction (see Appendix D.7). We also evaluate AutoMod at different filter levels ( $\alpha$ ) and find that for  $\alpha > 0$ , there is minimal gain in performance (Appendix D.5).

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

624

574 575

# 5 Case Studies: Implicit Hate, Context and Robustness

576 This section explores three case studies that exam-577 ine AutoMod's versatility and robustness.

Counterfactual Analysis: Unlike explicit hate speech, implicit hate is harder for systems to detect due to the usage of stereotypes, insinuations, or coded language – which may evade traditional 581 582 keyword-based filters. Therefore, detecting such content requires a deep understanding of context, 583 intent, and the subtle ways in which language can 584 perpetuate harm. For this study, we select 110 false negatives from SBIC and manually review them 586 to ensure that they are devoid of any slurs, while still remaining offensive. We then programmati-588 cally replace terms corresponding to demographic groups with offensive slurs associated with the same groups (e.g. replacing the phrase "Black People" with the *n*-word). When these counterfactual examples are passed to AutoMod (with  $C_A = C$ ), we observe 100% recall. Some examples from our 594 counterfactual set are presented in Table 6. This study highlights AutoMod's heavy bias towards us-596 ing slurs as an indication of hate and poor under-598 standing of context which allows implicit hate to easily evade detection.

Policy Adherence Evaluation: Twitch's Community Guidelines explicitly recognize the importance of context in content moderation: "At Twitch, we allow certain words or terms, which might otherwise violate our policy, to be used in an empow-604 ering way or as terms of endearment when such intent is clear." To study AutoMod's adherence to this policy, we manually select 20 non-slur sensitive fragments from SBIC and prompt the GPT-40 model to generate statements that incorporate these fragments in a non-offensive manner including us-610 age of the term in an educational or empowering 611 context. The prompt used for generation and some 612 examples are provided in Appendix E.2. We test these examples (with  $C_A$  set to C) with the filtering 614 level varying between Some Filtering ( $\alpha = 2$ ) and 615 616 Maximum Filtering ( $\alpha = 4$ ) for each filter. Under the former, AutoMod flags 89.5% of examples 617 while in the latter setting, it flags 98.5% examples. 618 This behavior is in contrast to the aforementioned policy. While a streamer can manually configure 621 AutoMod to allow these sensitive fragments when a stream is likely to elicit their usage in non-offensive contexts, this could give a free pass to bad actors

to perpetuate hate. This study further underscores the need for context-aware automated moderation.

**Robustness:** We measure AutoMod's robustness to minor, semantic-preserving perturbations of input text, given its seeming reliance on direct slurs. For this study, we prompt GPT-40 to make subtle changes (such as adding spaces/punctuation, altering spelling etc.) to 50 manually chosen sensitive fragments from the SBIC dataset. We employ 6 perturbation methods (listed in Table 7). For each fragment we also generate one example that uses the fragment as is (Unperturbed). The prompt used for generation and the definitions of the perturbation methods are provided in Appendix E.3. We observe that the recall drops from 100% to 4%on some of our perturbations. This indicates that malicious actors can easily evade AutoMod. However, it is worth noting that AutoMod was able to detect all perturbations (except the reverse one) of certain frequent terms such as the n-word. This indicates some degree of robustness for highly sensitive terms.

#### 6 Discussion and Future Work

AutoMod exhibits poor recall and F1-scores across datasets, compared to SoTA language models (Figure 2), and even GPT-2 (Sap et al., 2020). Our analysis of the false negatives and positives (§4.2 and Appendix D.3), along with the case studies  $(\S5)$  collectively explain how the poor performance of AutoMod arises from a lack of contextual understanding of hate speech. Regardless, Twitch's effort in the development of AutoMod is commendable for the flexibility it provides, with most other platforms lacking such tools. Our audit serves to highlight current gaps in AutoMod's capabilities and we make a good faith recommendation for Twitch to address these issues for greater user-safety. There are several directions for future work, the most of direct of which is to evaluate Twitch's Smart Detection feature for text, and expand the audit to audiovisual content. Considering other platforms and languages, particularly given the availability of multilingual hate speech data (Arora et al., 2023a) is of critical importance. More nuanced audits for black-box systems may be possible by leveraging techniques such as model reconstruction attacks (Tramèr et al., 2016) for reverse engineering. In conclusion, we hope our study serves as a blueprint for further audits of decision-making systems operating in complex socio-technical environments.

#### Limitations

674

While we attempt to be comprehensive in our evaluation of hate speech moderation on Twitch, several 676 key limitations remain. First, we only consider a 677 single platform in our study, and do not compare AutoMod to any other deployed text moderation. 679 As of he beginning of this study, we could not find any other suitable candidates. However, alternatives such as Mistral's Moderation tool Mistral AI Team (2024) have since emerged which would make for an interesting follow-up study, particularly in the context of other studies of open-source LLMs for moderation (Kumar et al., 2024). Crossplatform studies could offer insights on whether longstanding traditions of moderation on "longform," non-live textual platforms have any bearing 689 on the emergent moderation practices on platforms like Twitch, and vice versa. Applying our method-691 ology to other platforms could evolve it into a comprehensive audit blueprint capable of assessing various message-based moderation systems for their 695 handling of hateful content. Second, we largely investigate both hateful and benign messages in 697 a non-conversational context, *i.e.* messages are passed one-by-one without any simulation of dialogue. More work is needed to curate datasets of dialogue that contain both hateful and benign con-701 tent. In addition, investigating in-group/out-group dynamics may uncover differential treatment and moderation of content based on group affiliations. 703 In the case of AutoMod, however, we speculate that our results would have remained unchanged as the 705 policy violation case studies already demonstrate a 706 lack of contextual awareness. Third, our study is focused entirely on text in English. Considering other languages and modalities such as embedded text in images would be a critical direction for future 710 work. 711

#### 712 Ethical Considerations Statement

Exposure to Offensive Content: In conducting 713 this research, our team dealt with a significant 714 amount of offensive content. All authors were 715 aware of the nature of the work and consented to 716 view such content. It is important to note that no 717 individuals apart from the authors were exposed to 718 this material given our use of siloed experimental 719 720 streams. Moreover, the harmful content used in this work was sourced from open-source datasets. 721

**Legal Compliance:** While conducting this research, we did post content that violates Twitch's terms of service. However, this action was taken within the legal boundaries established by the *Sandvig v. Barr* case (Gilens, Naomi and Williams, Jamie, 2020), which protects such research activity.

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

**Potential for Adverse Impact:** We demonstrated how perturbations to offensive text circumvent Twitch's existing moderation systems, and there is a potential adverse impact wherein Twitch users might exploit these examples. However, these techniques and attacks are already welldocumented and have been extensively explored in the literature concerning NLP classifiers. If a hateful actor is driven to spread hate on Twitch, it is unlikely that they are unaware of these rather unsophisticated perturbations.

**Impact on Twitch Developers:** We also recognize the unfortunate, unlikely possibility that our study inadvertently increased the burden for Twitch's internal moderation developers. Twitch has added measures to protect its users, and we do not wish to undermine these efforts. Instead, we intend to highlight existing failures while advocating for careful audits of deployed systems. We will contact Twitch with our findings prior to the publication of this work, ensuring that they are aware of their system vulnerabilities and can take appropriate action.

#### References

- 47 U.S.C. § 230. 1996. Protection for private blocking and screening of offensive material. Available at: https://www.law.cornell.edu/uscode/ text/47/230.
- Arnav Arora, Maha Jinadoss, Cheshta Arora, Denny George, Haseena Dawood Khan, Kirti Rawat, Seema Mathur, Shivani Yadav, Shehla Rashid Shora, Rie Raut, et al. 2023a. The uli dataset: An exercise in experience led annotation of ogbv. *arXiv preprint arXiv:2311.09086*.
- Arnav Arora, Preslav Nakov, Momchil Hardalov, Sheikh Muhammad Sarwar, Vibha Nayak, Yoan Dinkov, Dimitrina Zlatkova, Kyle Dent, Ameya Bhatawdekar, Guillaume Bouchard, and Isabelle Augenstein. 2023b. Detecting harmful content on online platforms: What platforms need vs. where research efforts go. *Preprint*, arXiv:2103.00153.
- Bundesamt für Justiz. 2022. Gesetz zur Verbesserung der Rechtsdurchsetzung in sozialen Netzwerken (Netzwerkdurchsetzungsgesetz – NetzDG).

- 773 774
- 788
- 790 791
- 792
- 794
- 797
- 801

805

807

811

- 816 817
- 818
- 819 820
- 821 822

823

824 825

- Jie Cai, Aashka Patel, Azadeh Naderi, and Donghee Yvette Wohn. 2024. Content moderation justice and fairness on social media: Comparisons across different contexts and platforms. In Extended Abstracts of the CHI Conference on Human Factors in Computing Systems, CHI '24, page 1–9. ACM.
- Sarah H. Cen and Rohan Alur. 2024. From transparency to accountability and back: A discussion of access and evidence in ai auditing. Preprint, arXiv:2410.04772.
- Stevie Chancellor, Jessica Annette Pater, Trustin Clear, Eric Gilbert, and Munmun De Choudhury. 2016. #thyghgapp: Instagram content moderation and lexical variation in pro-eating disorder communities. In Proceedings of CSCW, CSCW '16, pages 1201-1213, New York, NY, USA. ACM.
- Eshwar Chandrasekharan, Umashanthi Pavalanathan, Anirudh Srinivasan, Adam Glynn, Jacob Eisenstein, and Eric Gilbert. 2017. You can't stay here: The efficacy of reddit's 2015 ban examined through hate speech. In Proceedings of ACM Human Computer Interaction, 1(Computer-Supported Cooperative Work and Social Computing):31:1-31:22.
- Eshwar Chandrasekharan, Mattia Samory, Shagun Jhaver, Hunter Charvat, Amy Bruckman, Cliff Lampe, Jacob Eisenstein, and Eric Gilbert. 2018. The Internet's Hidden Rules: An Empirical Study of Reddit Norm Violations at Micro, Meso, and Macro Scales. Proceedings of the ACM on Human-Computer Interaction, 2(CSCW):1-25.
- Jonathan P. Chang and Cristian Danescu-Niculescu-Mizil. 2019. Trajectories of blocked community members: Redemption, recidivism and departure. In Proceedings of WWW, pages 184–195.
- Le Chen, Ruijun Ma, Anikó Hannák, and Christo Wilson. 2018a. Investigating the impact of gender on rank in resume search engines. In Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems, CHI '18, page 1-14, New York, NY, USA. Association for Computing Machinery.
- Le Chen, Ruijun Ma, Anikó Hannák, and Christo Wilson. 2018b. Investigating the impact of gender on rank in resume search engines. In Proceedings of the 2018 chi conference on human factors in computing systems, pages 1–14.
- Le Chen, Alan Mislove, and Christo Wilson. 2015. Peeking beneath the hood of uber. In Proceedings of the 2015 internet measurement conference, pages 495-508.
- Benjamin Edelman and Michael Luca. 2014. Digital discrimination: The case of airbnb.com. SSRN Electronic Journal.

Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Divi Yang. 2021. Latent hatred: A benchmark for understanding implicit hate speech. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 345–363, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

- EU. 2000. Articles 12-14, european union e-commerece directive.
- Casey Fiesler, Jialun Jiang, Joshua McCann, Kyle Frye, and Jed Brubaker. 2018. Reddit rules! Characterizing an Ecosystem of Governance. In Proceedings of the International AAAI Conference on Web and Social Media, volume 12.
- Mirko Franco, Ombretta Gaggi, and Claudio E. Palazzi. 2023a. Analyzing the use of large language models for content moderation with chatgpt examples. In Proceedings of the 3rd International Workshop on Open Challenges in Online Social Networks, OASIS '23, page 1-8, New York, NY, USA. Association for Computing Machinery.
- Mirko Franco, Ombretta Gaggi, and Claudio E. Palazzi. 2023b. Integrating Content Moderation Systems with Large Language Models. ACM Trans. Web.
- S.M. Gaddis. 2018. Audit Studies: Behind the Scenes with Theory, Method, and Nuance. Methodos Series. Springer International Publishing.
- Gilens, Naomi and Williams, Jamie. 2020. Federal Judge Rules It Is Not a Crime to Violate a Website's Terms of Service. [Online; accessed 19. Jan. 2025].
- Tarleton Gillespie. 2018. Custodians of the Internet: Platforms, Content Moderation, and the Hidden Decisions that Shape Social Media. Yale University Press.
- Robert Gorwa, Reuben Binns, and Christian Katzenbach. 2020. Algorithmic content moderation: Technical and political challenges in the automation of platform governance. Big Data & Society, 7(1):2053951719897945.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, and et al. 2024. The llama 3 herd of models. Preprint, arXiv:2407.21783.
- Joanne E Gray and Nicolas P Suzor. 2020. Playing with machines: Using machine learning to understand automated copyright enforcement at scale. Big Data & Society, 7(1):2053951720919963.
- Catherine Han, Joseph Seering, Deepak Kumar, Jeffrey T. Hancock, and Zakir Durumeric. 2023. Hate raids on twitch: Echoes of the past, new modalities, and implications for platform governance. Preprint, arXiv:2301.03946.

Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi,

Maarten Sap, Dipankar Ray, and Ece Kamar. 2022.

ToxiGen: A large-scale machine-generated dataset

for adversarial and implicit hate speech detection.

In Proceedings of the 60th Annual Meeting of the

Association for Computational Linguistics (Volume

1: Long Papers), pages 3309–3326, Dublin, Ireland.

Jacob Foster. 2023. tmi.js. [Online; accessed 20. Jan.

Brian Keegan and Casey Fiesler. 2017. The evolution and consequences of peer producing wikipedia's

Daphne Keller, Paddy Leerssen, et al. 2020. Facts and

where to find them: Empirical research on internet

platforms and content moderation. Social media and

democracy: The state of the field and prospects for

Mahi Kolla, Siddharth Salunkhe, Eshwar Chandrasekha-

ran, and Koustuv Saha. 2024a. Llm-mod: Can large

language models assist content moderation? In Ex-

tended Abstracts of the CHI Conference on Human

Factors in Computing Systems, CHI EA '24, New York, NY, USA. Association for Computing Machin-

Mahi Kolla, Siddharth Salunkhe, Eshwar Chandrasekha-

ran, and Koustuv Saha. 2024b. LLM-Mod: Can

Large Language Models Assist Content Moderation?

CHI EA '24: Extended Abstracts of the CHI Confer-

ence on Human Factors in Computing Systems, pages

Jeff Kosseff. 2019. The twenty-six words that created

Deepak Kumar, Yousef AbuHashem, and Zakir Du-

Yao Lyu, Jie Cai, Anisa Callis, Kelley Cotter, and

John M. Carroll. 2024. "i got flagged for supposed

bullying, even though it was in response to someone harassing me about my disability.": A study of blind

tiktokers' content moderation experiences. In Pro-

ceedings of the CHI Conference on Human Factors

in Computing Systems, CHI '24, page 1-15. ACM.

J. Nathan Matias, Austin Hounsel, and Nick Feamster.

vertising policies. Preprint, arXiv:2103.00064.

port. [Online; accessed 19. Jan. 2025].

Meta. 2024. Community Standards Enforcement Re-

Danaë Metaxa, Michelle A Gan, Su Goh, Jeff Hancock,

and James A Landay. 2021a. An image of society:

Gender and racial representation and impact in image

search results for occupations. Proceedings of the

ACM on Human-Computer Interaction, 5(CSCW1):1-

2021. Software-supported audits of decision-making

systems: Testing google and facebook's political ad-

rumeric. 2024. Watch your language: Investigat-

ing content moderation with large language models.

the Internet. Cornell University Press.

Preprint, arXiv:2309.14517.

Association for Computational Linguistics.

rules. In Proceedings of ICWSM.

reform, pages 220-251.

- 890

2025].

erv.

1 - 8.

23.

- 891
- 896

898

900

901

905 906

907

911

914

908 909 910

- 912 913
- 915 916 917
- 921 922

- 925 926 927
- 928
- 930

931 932

933 935

936

Danaë Metaxa, Joon Sung Park, Ronald E. Robertson, Karrie Karahalios, Christo Wilson, Jeff Hancock, and Christian Sandvig. 2021b. Auditing algorithms: Understanding algorithmic systems from the outside in. Foundations and Trends® in Human–Computer Interaction, 14(4):272–344.

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

- Jakub Mikians, László Gyarmati, Vijay Erramilli, and Nikolaos Laoutaris. 2012. Detecting price and search discrimination on the internet. In Proceedings of the 11th ACM workshop on hot topics in networks, pages 79–84.
- Mistral AI Team. 2024. Mistral Moderation API. [Online; accessed 16 Feb. 2025].
- Annabelle Nyst. 2023. 134 Social Media Statistics You Need To Know For 2023 - searchenginejourhttps://www.searchenginejournal. nal.com. com/social-media-statistics/480507/. [Accessed 17-01-2025].
- Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. 2019. Dissecting racial bias in an algorithm used to manage the health of populations. Science, 366(6464):447-453.
- Victor Le Pochat, Laura Edelson, Tom Van Goethem, Wouter Joosen, Damon McCoy, and Tobias Lauinger. 2022. An audit of facebook's political ad policy enforcement. In 31st USENIX Security Symposium (USENIX Security 22), pages 607-624, Boston, MA. USENIX Association.
- Manoel Horta Ribeiro, Justin Cheng, and Robert West. 2023. Automated content moderation increases adherence to community guidelines. Preprint, arXiv:2210.10454.
- Nadia Ruiz Bravo and Maryam Roshan. 2022. The political turn of twitch - understanding live chat as an emergent political space.
- Christian Sandvig, Kevin Hamilton, Karrie Karahalios, and Cédric Langbort. 2014. Auditing algorithms : Research methods for detecting discrimination on internet platforms.
- Maarten Sap, Saadia Gabriel, Lianhui Oin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 5477–5490, Online. Association for Computational Linguistics.
- Brennan Schaffner, Arjun Nitin Bhagoji, Siyuan Cheng, Jacqueline Mei, Jay L Shen, Grace Wang, Marshini Chetty, Nick Feamster, Genevieve Lakier, and Chenhao Tan. 2024. "community guidelines make this the best party on the internet": An in-depth study of online platforms' content moderation policies. In Proceedings of the CHI Conference on Human Factors in Computing Systems, CHI '24, page 1-16. ACM.

Joseph Seering and Sanjay R. Kairam. 2022. Who moderates on twitch and what do they do? quantifying practices in community moderation on twitch. *Proc. ACM Hum.-Comput. Interact.*, 7(GROUP).

991

992

994

997

999 1000

1001

1002

1003 1004

1005

1006

1007

1008

1009

1010

1011

1012 1013

1014

1015

1016 1017

1018

1019 1020

1021

1023

1024

1025

1026

1027 1028

1029

1030

1031

1032

1033

1034

- Till Speicher, Muhammad Ali, Giridhari Venkatadri, Filipe Nunes Ribeiro, George Arvanitakis, Fabrício Benevenuto, Krishna P. Gummadi, Patrick Loiseau, and Alan Mislove. 2018. Potential for discrimination in online targeted advertising. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pages 5–19. PMLR.
- Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. 2016. Stealing machine learning models via prediction {APIs}. In 25th USENIX security symposium (USENIX Security 16), pages 601–618.
- Twitch. 2024a. How to Use AutoMod. [Online; accessed 19. Jan. 2025].
- Twitch. 2024b. Transparency Reports. [Online; accessed 18. Jan. 2025].
  - Twitch. 2025. Chat & Chatbots. [Online; accessed 19. Jan. 2025].
- TwitchTracker. 2025. Twitch channels, games and global statistics.
- Bertie Vidgen, Kayla Chatelon, Scott Hale, Helen Margetts, and Dong Nguyen. 2021a. Learning from the worst: Dynamically generated datasets to improve online hate detection. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4216–4231. Association for Computational Linguistics.
- Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2021b. 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.
- D Wakabayashi. 2019. Legal shield for websites rattles under onslaught of hate speech. *The New York Times*, *Aug*, 6.
- X. 2024. Dsa transparency report october 2024.

- 1038
- 1039 1040
- 1042
- 1043
- 1044
- 1045
- 1046 1047
- 1048

- 1052
- 1053 1054
- 1055
- 1056
- 1057 1058
- 1059
- 1060
- 1061 1062
- 1063
- 1064

1065

1066

1068

- 1069
- 1070
- 1071

1072

1074

1077

This Appendix provides additional information corresponding to each section of the main paper as follows:

- 1. Platform Choice and Twitch Moderation Details (Appendix A): Provides additional details about Twitch Moderation options referenced in §3.1.
- 2. Dataset Details (Appendix B): Detailed descriptions of the 4 datasets of hate speech content introduced in  $\S3.2$ .
- 3. Experimental Setup Details (Appendix C): Breakdown of the steps followed to send and receive messages from Twitch, expanding upon §4.1.
- 4. Further Results (Appendix D): Building on results in §4.2 with analysis of false negatives and false positives from AutoMod, SBIC data threshold, different AutoMod filter levels, prefiltering bias, and quality control analysis of filter-specific datasets.
- 5. Ablation Details (Appendix E): Information on the methods used to construct counterfactuals, policy adherence samples and perturbed examples in §5.
- 6. Filter- and Community-Specific Subset Extraction (Appendix F): Dataset-wise breakdown of the procedure followed to obtain subsets for experiments in  $\S4.2$ .

#### **Platform Choice and Twitch** Α **Moderation Details**

As an extension of the Platform Exploration and Twitch Moderation section  $(\S3.1)$ , we provide more details about the investigation guiding our platform choice, the Twitch interface, and its internal documentation.

# A.1 Platform Survey

In Table 3, we provide details about the 3 platforms that we ended up choosing from in terms of their 1073 suitability for the audit. While both Reddit and Discord offer siloed environments for experimentation, 1075 neither provides native machine learning models 1076 for offensive text detection.



(a) Twitch Moderation Pyramid What AutoMod catches at each level

Level	Description
Level 0	No filtering.
	Tip: Words and terms you have added to your channel's blocked terms will still be blocked even if AutoMod is turned off. This can prevent messages from being sent.
Level 1	Some filtering on discrimination, and Smart Detection only.
Level 2	Some filtering on discrimination, sexual content, and Smart Detection, more filtering on hostility.
Level 3	More filtering on discrimination, sexual content, hostility, and Smart Detection.
Level 4	More filtering on discrimination, sexual content, profanity and Smart Detection, and most filtering on hostility.

(b) AutoMod Moderation Level Distinctions from Level 0 to 4

Figure 4: Official Twitch Documentation of Moderation Tools: (a) Moderation tools provided for each category of user and platform-wide (b) Twitch Documentation on AutoMod Moderation Levels from 0 to 4.

# A.2 Twitch Automated Moderation

In Figure 4a, we show the Twitch Moderation Pyramid (Twitch, 2025). As stated in §3.1, we focus on viewer-to-streamer and streamer-to-streamer which are moderated via AutoMod and Chat Filters. In Figure 4b, we show the various levels Twitch offers moderators to customize filtering strength.

1078

1079

1080

1081

1083

1084

1085

1086

1088

1089

1090

1091

1092

1094

1095

1096

1097

1098

Table 4 describes Twitch's Official Content Categories and their differentiating definitions as of 04 July, 2024. The audit's investigation focuses on the effectiveness of AutoMod overall performance with respect to the Discrimination and Slurs content category. More specifically, the datasets used provide ground truth relating to the subsections of Discrimination and Slurs as seen in Table 4, Disability, SSG, Misogyny, and RER.

The dashboard is customizable allowing the streamer to view in real-time their stream's AutoMod queue and to manually approve or deny detected messages as shown in Figure 5. The AutoMod queue also displays the associated content

Platform	<b>Closed environment</b>	Moderation tools
Discord	Private server with permissions for everyone but admin disabled	Image hashing and "ML powered tech" for child sexual abuse material; No native ML models for text, relies on keyword detection
Reddit	Via private subreddits	Needs plugins for ML models and subreddit specific rules for moderation; Well-investigated (Kolla et al., 2024b; Ku- mar et al., 2024; Franco et al., 2023b)
Twitch	Untagged streams are private	Customisable moderation levels using native ML models; Also has keyword lists and Smart Detection to train model on the actions of human moderators

Table 3: Comparing candidate platforms for the audit study

category so as to inform the streamer for Auto-Mod's reasoning behind moderation. Highlighted segments also help the user focus on the problematic fragments of moderated messages.

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

In evaluating the overall performance of Auto-Mod on moderating various categories of hateful content, mere binary evaluations were deemed primary but not comprehensive to effectively determine performance in terms of accuracy and precision. Hence, AutoMod's moderation reasons with respect to Twitch's pre-defined categories of moderation provided more insight. To programmatically access the streamer's AutoMod queue, Twitch Pubsub event subscriptions were used to grab API results whenever AutoMod logged a message in the queue for moderation. However, we found that AutoMod's internal categories did not match Twitch Documentation and Policies on Content Moderation across the platform. Figure 6 depicts a snapshot of Twitch API result from automodqueue event subscription detailing the message and its metadata. This metadata include topics and category of the moderated message from which we map the categories semantically.

#### **B** Dataset Details

All datasets used are open-source and permitted for research. All of our data is in English.

## B.1 Dynahate

1127The Dynahate dataset (Vidgen et al., 2021b) con-1128tains approximately 41,000 entries, each labeled as1129either 0 (Not hateful) or 1 (Hateful) and is gener-1130ated and labelled by trained annotators over four1131rounds of dynamic data creation. The dataset in-1132cludes 54% of hateful examples. In our work, we1133use three columns from the Dynahate dataset: *text*,

*label*, and *target*. The 'target' column specifies the community toward which the hate is directed, making it particularly useful for dividing the data into filter and community-specific subsets. 1134

1135

1136

1137

1138

1163

# **B.2 SBIC**

The Social Bias Inference Corpus (SBIC) (Sap 1139 et al., 2020) is a large-scale dataset with over 150k 1140 annotations across 34,000 social media posts, cap-1141 turing nuanced biases and stereotypes in online 1142 language. Its focus on offensive content, implied 1143 stereotypes, and target demographics makes it par-1144 ticularly useful for auditing Twitch's moderation 1145 system, especially in areas related to group-based 1146 bias and contextual offensiveness. We sample 20k 1147 examples from the SBIC training data for our work. 1148 Since SBIC provides an average annotator offen-1149 siveness score, we set a threshold of 1 (all annota-1150 tors agreeing on the offensiveness of an example) 1151 to evaluate the overall effectiveness of AutoMod, 1152 resulting in a hate-to-non-hate ratio of 1:3. For 1153 further division into hateful subsets related to spe-1154 cific filters, we use a threshold of 0.5, yielding a 1155 total of 8,748 examples. The SBIC dataset con-1156 tains numerous features, but for our analysis, we 1157 utilize the following: post, targetMinority, target-1158 Category, offensiveYN (indicating whether the ex-1159 ample is offensive, non-offensive, or ambiguous), 1160 and annotator-related columns for aggregating the 1161 offensive score. 1162

# B.3 ToxiGen

The ToxiGen dataset (Hartvigsen et al., 2022) is a1164large-scale dataset containing approximately 274k1165toxic and benign statements about 13 minority1166groups. This dataset is particularly useful for our1167work as it consists of synthetically generated im-1168

Moderation Category	Explanation of Hate Speech Category
Sexual Content	Words or phrases referring to sexual acts and/or anatomy.
Discrimination and Slurs	Includes race, religion, gender-based discrimination. Hate speech falls under this category.
I) Disability	Demonstrating hatred or prejudice based on perceived or actual mental or physical abilities.
II) SSG	Demonstrating hatred or prejudice based on sexual identity, sexual orien- tation, gender identity, or gender expression.
III)Misogyny	Demonstrating hatred or prejudice against women, including sexual objec- tification.
IV) RER	Demonstrating hatred or prejudice based on race, ethnicity, or religion.
Hostility	Provocation and bullying, sexual harassment.
Profanity	Expletives, curse words, and vulgarity. This filter especially helps those who wish to keep their community family-friendly.
Smart Detection	Detects unwanted messages (including spam) based on moderation actions taken in your channel.

Table 4: AutoMod Content Categories and Subcategory details

plicit hate speech, allowing us to evaluate the ef-1169 fectiveness of filters on implicit hate content. The 1170 authors use the GPT-3 model to generate implicit 1171 hate examples by providing human-curated toxic 1172 and non-toxic examples as prompts. For detailed 1173 information on the data generation methods, we 1174 refer the reader to the ToxiGen paper. ToxiGen 1175 includes various features, but we use the following 1176 in our work: prompt, generation (text generated by 1177 the model—examples we use for auditing), prompt 1178 label, and roberta prediction for the generated sen-1179 tence. For our experiment, we randomly sample 1180 20k hateful examples from the dataset, selecting 1181 entries with toxic prompts and a roberta\_prediction 1182 score between 0.8 and 1.0. Similarly, we sample 1183 20k non-hateful examples with non-toxic prompts 1184 and a *roberta\_prediction* score between 0 and 0.2. 1185

#### B.4 Implicit Hate Corpus(IHC)

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

The Implicit Hate Corpus (IHC) (ElSherief et al., 2021) is a dataset containing approximately 6,000 examples of implicit hate. It includes both the target and implied meaning of hateful statements, collected from hate communities and their followers on Twitter. In our work, we use this dataset to evaluate the effectiveness of all discrimination filters. Since the dataset is not large enough to be divided into filter-specific subsets, we only remove examples that are not aligned with discrimination.

From this dataset, we utilize the *post* and *target* columns.

1197

1198

1199

1200

1201

1203

1204

1205

1206

1207

1208

1209

#### **C** Experimental Setup Details

Here we describe in detail the setup and components needed to build the experimental pipeline (Figure 1) that was briefly described in §4.1. The critical steps for implementation include: 1) bot creation to send messages, including adherence to security measures, registration, and scope handles, 2) AutoMod event subscriptions for sorting messages into moderated and unmoderated, and 3) results handling.

#### C.1 Sending Messages Programatically

The following steps are needed to *send* messages 1210 programatically at scale: Step 1: Create a chatbot-1211 In order to handle the creation of chatbots, Twitch 1212 requires a two-factor authenticated phone number 1213 and email address for each account. The bots cre-1214 ated for implementation utilize free online phone 1215 numbers for 2FA, Burner phone numbers (Veri-1216 fied Phone Numbers only), and personal phone 1217 numbers. The different bots also require an ap-1218 plication registration in the associated verified-1219 account's Twitch Developer Portal. This process 1220 required the specification of a unique chat appli-1221 cation name, OAuth Redirect URL specification: 1222 https://localhost:3000, and the bot's functionality: 1223



Figure 5: Moderation View from Streamer's Account

Chat Bot. This application enables Twitch Devel-1224 1225 opers to extract application-specific Client-ID, and Client-Secrets; Step 2: Authenticate and register 1226 the chatbot- Twitch enables authentication under 1227 OAuth2.0 using Twitch OIDC Authorization Code 1228 Grant Flow in order to grant the application specific 1229 access to Twitch HTTPS resources. This access, for 1230 example, is specified as chat:read, chat:write 1231 access for IRC Chatbots. Once the Au-1232 1233 thorization code from HTTPs POST requests to https://id.twitch.tv/oauth2/token is re-1234 ceived, the app access token and refresh tokens 1235 are retreived and used for API calls. As each application's access token used to connect to Twitch 1237 Chat and AutoMod queues, we automatically han-1238 dle the renewal of these authorized tokens using 1239 the application's refresh token and reconnect ac-1240 cordingly when its expiration time has reached; 1241 Step 3: Send messages at scale with appropri-1242 ate timing configuration- We create one instance 1243 of the experimental pipeline by including a mes-1244 senger bot, channel bot, and Pubsub bot (three-bot 1245 1246 configuration). The messenger bot which sends up to 5 messages with a 4 second wait between 1247 each message. At every iteration, we pause for 3.5seconds due to Twitch's Chat Rate Limits (Twitch, 1249 2025) allowing for less than 20 messages per 30 1250 seconds for normal chat bot accounts. This ensures 1251

that each message is sent and processed into the stream chat without duplication or misses due to connection delays. We conducted iterative tests concluding that this configuration of pauses and message counts allowed for prolonged experiment run times due to the large number of messages from our four datasets.

1252

1253

1254

1255

1256

1257

1258

1259

1260

# C.2 Receiving moderated and unmoderated messages

To receive and sort the messages we need to cre-1261 ate additional bots and parse their outputs. This 1262 connection is facilitated by the use of tmi.js 1263 JavaScript package that creates a connection to 1264 the Twitch IRC server using tmi.js servers (Jacob 1265 Foster, 2023). Twitch IRC presents a reduced-1266 functionality RFC1459 and IRCv3 Message Tag 1267 specification (Twitch, 2025) for parsing messages 1268 to and from a specified Twitch channel. This extrac-1269 tion of non-moderated messages is handled by the 1270 Receiver Bot. The Pubsub bot uses PubSub event 1271 subscription to extract messages from the AutoMod 1272 queue in real time which are accessible from the 1273 channel's creator dashboard. Using twitchAPI.js, 1274 we enable a websocket connection to Twitch ma-1275 chines providing Pubsub services from which regis-1276 tered chat bots listen for AutoMod queue items and 1277 receive message metadata in a websocket frame. 1278



Figure 6: API Call result of Moderated result mapped to Content Category

#### C.3 Parsing results

1280

1281

1282

1283

1284

1285

1288

1289

1290

1291

1292

1293

1294

1295

1297

1298

1299

1300

1303

1304

1305

1306

1308

The json received from Pubsub provides information about the moderated fragment, and an internal label for moderated content such as Ableism, Misogyny, Racism, and Homophobia (Figure 6). While these categories are not detailed in AutoMod documentation, we infer an injective relationship between the content categories in AutoMod documentation and internal categories from the API call. With these labels, we are able to investigate AutoMod's stated reasons for moderation and conduct further category-specific analysis in §4.2.

#### C.4 Challenges

Despite being under the rate limit, we suspect that the increase in fraudulent activity across Twitch as reported in 2024's Transparency Report (Twitch, 2024b) has led to the inaccurate but frequent permanent banning of our chat bots. As a result, this lead to the creation of at least 30 different accounts. After conducting several initial experiments and analyzing the results, we discovered that certain messages never appeared in the receiver bot's inbox and, consequently, did not trigger any Pub/Sub events. To investigate further, we manually passed some of these messages and found that Twitch had prefiltered them, preventing them from being sent. To address this, we enhanced our code to log all messages-both those moderated by AutoMod and those that were not-into a CSV file. By comparing this log against the original input data and

filtering out the missing messages, we successfully1309identified the set of prefiltered messages.1310

#### **D** Further Results

All our experiments with language models were carried out on a workstation with 6 A6000 GPUs.

#### **D.1** Evaluating the Chat Filter

The viewer-controlled chat filter can only act at the 1315 granularity of the broad Discrimination and Slurs 1316 category, and does not trigger events for the Pub-1317 sub bot we use to detect moderation events which 1318 makes experiments at scale challenging as man-1319 ual inspection is required to identify moderation, 1320 indicated by a "\*\*\*" (star pattern). To audit the 1321 alignment between AutoMod and the Chat Filter, we 1322 randomly sample 200 examples from SBIC consisting of equally many moderated and non-moderated 1324 examples and pass them through the chat filter. We 1325 observe that the Chat Filter is unable to moder-1326 ate any of the non-moderated examples (Figure 7). 1327 We infer that the Chat Filter behaves similarly to AutoMod at  $\alpha = 4$ . 1329

#### **D.2** Filter Precision

The filter precision for filter  $\mathcal{F}_i$ , denoted by  $P_{\mathcal{F}_i}$  is defined as:

$$P_{\mathcal{F}_i} = \frac{\sum\limits_{x \sim \mathcal{D}_{c_i}} \mathbb{1}(\mathcal{F}_i(x) = 1)}{\sum\limits_{j \in \widetilde{\mathcal{S}_c}} \left( \sum\limits_{x \sim \mathcal{D}_{c_j}} \mathbb{1}(\mathcal{F}_i(x) = 1) \right)}$$

1330

1311

1312

1313

freely and incapable of being civilized mollykim123: "her physical appearance makes her unattractive; she's not good looking at all ----- - Hey dude, \*\*\*\*\*\*\* - No way man"

mollykim123: wheelchair-bound \*\*\*\*\* in your class gives you dirty looks for fiddling with your leg brace in class

(a) Twitch Simulation: User with chat filters turned on for Disability, Misogyny, RER, and SSG excluding profanity

**mollykim123**: all lgbtg people should live freely

mollykim123: people of color are violent and incapable of being civilized

mollykim123: "her physical appearance makes her unattractive; she's not good looking at all ----- - Hey dude, youre gay - No way man"

mollykim123: wheelchair-bound bitch in your class gives you dirty looks for fiddling with your leg brace in class

(b) Twitch Simulation: User without chat filters on

Figure 7: Demonstrating behavior similar to AutoMod whereby false positives are prevalent and fails our target group analysis standards as seen in (a). Original Messages are provided in (b)

It is the number of examples adhering to criterion  $c_i$  (i.e., from  $\mathcal{D}_{c_i}$ ) that were moderated by filter  $\mathcal{F}_i$ upon the number of examples across all subsets that were moderated by filter  $\mathcal{F}_i$ . The precision of each filter is shown in Figure 8.

1331

1332

1333

1334

1335



Figure 8: The precision of filters in AutoMod. Precision measures the specificity of a filter.

#### **False Negative and False Positive Analysis D.3**

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1352

1353

#### D.3.1 False Negative Analysis

**Setup** For this experiment, we use three of our datasets - SBIC, DynaHate and ToxiGen. We randomly choose a dev set of 10,000 examples from these three datasets and use the rest as test data. We use the dev set to tune our prompt and develop our answer extractor. We prompt the llama-3.3-70B-instruct model (Grattafiori et al., 2024) to identify if the provided text contains swear word(s) and record the model outputs for the test data. During both dev and test phases, we only use the False Negatives.

**Prompt** The prompt used for this experiment was a zero-shot prompt that instructed the model to identify the presence of swear word(s) in the provided text and return a binary label. The exact prompt is given below:

System Prompt: "You will be given a 1354 text. Your task is to act like a human 1355 judge and identify if the text provided 1356 to you contains swear words. Your re-1357 sponse should be a dictionary with a 1358 single key - "has\_swear\_words" which 1359 corresponds to whether the provided text has swear words. You should only re-1361 spond in the format as described below. 1362 DO NOT RETURN ANYTHING ELSE. 1363 START YOUR RESPONSE WITH '{'. [response format]: 1365 { 1366 "has\_swear\_words": "True" if the 1367 provided text has swear words else "False" 1369 }" User Prompt: "The provided text is: [text]: {text}" 1372 1373

## **D.3.2** False Positives

AutoMod exhibits a low False Positive Rate (FPR) 1374 overall, and most of the false positives occur in the 1375 DynaHate dataset — on which AutoMod has a rel-1376 atively higher FPR. Similar to the False Negative 1377 analysis, we hypothesize this is due to a high number of swear words in DynaHate and AutoMod's 1379 overreliance on profanity as a signal for hate. We 1380 perform an experiment similar to the False Nega-1381 tive Analysis to detect the presence of swear words 1382 1383on DynaHate's false positive set. We find that1384nearly 73% of the false positives from DynaHate1385contain profanity.

1386

1387

1388 1389

1390

1391

1392

1393

1394

1395

1396

1397

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

**Setup** The setup remains the same as Appendix D.3.1, except that we use only those examples from the test set which are false positives from the DynaHate dataset.

#### D.4 Performance Metrics at different values of SBIC threshold

The SBIC dataset provides annotator agreement scores for the offensiveness of each sentence, which we used to classify examples: sentences with agreement scores below threshold were categorized as non-hateful, while those with scores of threshold or greater were classified as hateful. See Figure 9 for Moderation Rate(Recall) on different values of threshold for SBIC.



Figure 9: AutoMod's recall (with  $C_A = \tilde{C}$ ) vs. the offensiveness score threshold used for obtaining SBIC ground truth labels. Recall reduces when more nuanced examples (lower threshold) of hate are included.

#### D.5 Performance Metrics at different filter levels for IHC Dataset

Twitch offers five levels of filtration, with each filter adjustable from **No Filter** ( $\alpha = 0$ ) to **Maximum Filter** ( $\alpha = 4$ ). In our experiment, we apply the IHC dataset to all five filtering levels and measure recall for each, as shown in Table 5. The variation across the different levels is minimal, except for  $\alpha = 0$ , where AutoMod does not perform any moderation, and only the prefiltered text undergoes moderation.

Filter Level	Recall (%)
No filter	0.18
Less filter	5.31
Some filter	6.31
More filter	6.39
Maximum filter	6.42

Table 5: Recall values for IHC Dataset on various filterlevels supported by Twitch

### D.6 Experimental Setup for computing performance metrics of SoTA language models

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

We choose 7 different language models with parameters ranging from 2B to 70B to compare with AutoMod. We use the instruct variants of the models wherever available. We choose a random sample of 10,000 examples as a test set. We ensure that the test set is balanced between the two classes. The numbers shown for AutoMod in Figure 2 are also correspond to this test set. We prompt the language models with Twitch's Community guidelines and ask to identify if the provided comment violates the guidelines.<sup>4</sup> We use a zero-shot prompt for our experiment. We set temperature = 0 for reproducibility and use a top\_p value of 1. For models which support system prompts, we add the instruction and community guidelines in the system prompt and the comment to be labeled in the user prompt. For other models, we add all text in the user prompt itself. The prompt is shown below and the performance metrics of the different models alongside AutoMod are shown in Figure 2. To avoid redundancy, we only show the prompt in the System + User format.

#### D.6.1 Prompt

*System Prompt:* "You will be given a text. Your task is to act like an online human moderator on Twitch and judge if the provided comment conforms to the community guidelines of the platform (Twitch):

## Community Guidelines:

Twitch does not permit behavior that is motivated by hatred, prejudice or intolerance, including behavior that promotes or encourages discrimination, denigration, harassment, or violence based on the following protected characteristics: race, ethnicity, color, caste, national origin, immigration

<sup>&</sup>lt;sup>4</sup>For the prompt, we exclude part(s) of the Community Guidelines that are not relevant for text based moderation (e.g., guidelines for images/videos)

1448status, religion, sex, gender, gender identity, sexual1449orientation, disability, serious medical condition,1450and veteran status. We also provide certain pro-1451tections for age, which are expressly noted in the1452examples.

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1474

1475

1476

1477

1478

1479

1480

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

We define 'protected groups' as a subset of the population with a shared protected characteristic. Every Twitch user falls into multiple of these protected groups. Twitch affords every user globally equal protections under this policy.

For example, regardless of your intent, you may not:

- 1. Promote, glorify, threaten, or advocate violence, physical harm, or death against individual(s) or groups on the basis of a protected characteristic, including age.
- 2. Use hateful slurs, either untargeted or di-1464 rected towards another individual. We allow 1465 certain words or terms, which might otherwise violate our policy, to be used in an empow-1467 ering way or as terms of endearment when 1468 such intent is clear. We also make exceptions 1469 for slurs in music—and singing along to mu-1470 sic—as long as the song itself is not hateful 1471 and the slurs are not combined with other dis-1472 criminating or denigrating content. 1473
  - 3. Post a comment that dehumanizes or perpetuates negative stereotypes and/or memes.
    - 4. Post a comment that expresses inferiority based on a protected characteristic, for example, statements related to physical, mental, and moral deficiencies.
  - 5. Call for subjugation, segregation or exclusion, including political, economic, and social exclusion/segregation, based on a protected characteristic, including age. We do, however, allow discussions on certain topics such as immigration policy, voting rights for noncitizens, and professional sports participation as long as the content is not directly denigrating based on a protected characteristic.

6. Encourage or support the political or economic dominance of any race, ethnicity, or religious group, including support for white supremacist/nationalist ideologies. This does not include support for self-determination movements.

7. Expressions of contempt, hatred, or disgust 1495 based on a protected characteristic. 1496 8. Mock the event/victims or deny the occurrence 1497 of well-documented hate crimes, or deny the 1498 existence of documented acts of mass mur-1499 der/genocide against a protected group. 9. Make unfounded claims assigning blame to 1501 a protected group, or that otherwise intends 1502 to incite fear about a protected group as it 1503 relates to health and safety. 1504 10. Encourage the use of or generally endorsing 1505 sexual orientation conversion therapy. 1506 11. Support, promote, or be a member of a hate 1507 group, including sharing hate group propa-1508 ganda materials. 1509 The response should be a dictionary with a sin-1510 gle key - "label", which corresponds to whether 1511 the comment conforms to the community guidelines 1512 (nothate or hate). 1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1526

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

You should only respond in the format as described below. DO NOT RETURN ANYTHING ELSE. START YOUR RESPONSE WITH '{'.

[response format]:

{

"label": 'hate' if the comment violates community guidelines, else 'nothate'

- }"
  User Prompt: "The following is the given comment:
  - [text]: {comment}"

#### **D.7** Prefiltering Bias

The *n*-word is a highly offensive racial slur and is considered inappropriate in nearly all contexts barring anti-racist uses in, for example, art. The n-word's infamy may push system designers to prefilter the term as well as its derivative forms. We observe a high frequency of the *n*-word in the prefiltered examples as shown in Figure 10. While such system designs are done in good faith, they may induce biases in the system when the channel level moderation algorithms are not sophisticated enough and the blocklist used for pre-filtering is not exhaustive.<sup>5</sup> We hypothesize that due to the presence of the *n*-word, the system will perform much better at blocking hate directed towards Black

<sup>&</sup>lt;sup>5</sup>Here, by 'blocklist' we refer to the implicitly/explicitly defined list of words that the pre-filtering algorithm filters (Figure 10).

1540Folks than it would for the other communities. To1541test this hypothesis, we set  $C_A = \tilde{C} \setminus \{c_{RER}\}$  and1542record the outputs. We observe that on both SBIC1543and DynaHate combined, 37.4% of the recall on1544Black-related hate speech is due to prefiltering. For1545the Jewish and Muslim categories, this number is1546much lower (6.1% and 8.6% respectively).



Figure 10: Top 10 unigrams based on frequency from prefiltered examples after removing stop words. We observe that both variants of the n-word combined together appear close to 800 times.

#### D.8 Quality Control analysis for data subsets

1547

1548

1550

1552

1553

1554

1555

1556

1558

1559

1560

1562

1563

1564

1565

1568

1571

1572

1573

1575

In §4.2 we've created filter-specific subsets of data. It is however important to measure the specificity of these subsets as well to build confidence in our analysis. As long as the subsets are specific enough for their corresponding filters, we can have faith in our analysis. To measure this, we look at subsets from the filter point-of-view: For each subset, we turn on all filters except the filter that corresponds to that subset and record AutoMod's decisions. Then, we compute the recall for each subset. If a subset has too many examples that fit into criteria other than that of the filter the subset corresponds to, then such a subset should exhibit a higher recall when the relevant filter is turned off - as other filters would still be able to flag the examples with overlapping criteria. The subset-wise recall is shown in Figure 11. Except, SSG all other data subsets have a low recall when their corresponding filter is turned off. This implies that these subsets align strongly with the filter criterion, in the sense that they do not get moderated by filters other than the one they're supposed to test. For the SSG filter, the high recall is due to the pre-filtering. In both SBIC and DynaHate, we see that more than 50% of the examples from the SSG subset get pre-filtered. Since pre-filtering operates above the filter level, one can only attribute the non shaded portion of the each bar in Figure 11. With this justification,



Figure 11: Subset-wise recall when all filters are turned on except the one that shares the same category as the subset

we believe that our subsets are categorized well enough to be used as test sets for auditing the four filters. 1576

1577

1579

1580

1581

1582

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

#### E Ablation Details

#### E.1 Counterfactual Examples

For the Counterfactual study, we select 110 false negatives from SBIC and generated counter factual examples as per the process discussed in 5. Table 6 presents some demo examples from our counterfactual set.

#### E.2 Prompt for Policy Discrepancy Experiment

Policy Discrepancy Prompt to GPT4.0 model through ChatGPT (Interface): "Using the Twitch policy: 'At Twitch, we allow certain words or terms, which might otherwise violate our policy, to be used in an empowering way or as terms of endearment when such intent is clear. We also make exceptions for slurs in music—and singing along to music—as long as the song itself is not hateful and the slurs are not combined with other discriminating or denigrating content.'

Generate 10 examples using the phrase 'fragment' in positive or neutral contexts. Each example should:

- 1. Use the phrase in an empowering or endearing way 1602
- 2. Use the phrase in educational contexts to raise awareness 1603
- 3. Use the phrase in neutral or harmless contexts,<br/>aligning with Twitch's policy16051606

#### Content Warning: The table below contains bigotry that may be distressing to some readers.

Original (Not Moderated)	Counterfactual (Moderated)		
why are all <i>black people</i> fast. because all the slow ones are in jail.	why are all <i>n</i> fast. because all the slow ones are in jail.		
what do you call a <i>muslim</i> abortion bomb de- fusal	what do you call a <i>R</i> abortion bomb defusal		
why do <i>jewish</i> men get circumcised? Because jewish women won't touch anything that isn't 30% off.	why do $k$ men get circumcised? because $k$ women won't touch anything that isn't 30% off.		

Table 6: Comparison of original and counterfactual Texts: demonstrating how explicit slurs trigger moderation while implicit hate speech remains unflagged

Perturbation	Example	Moderation Rate (%)
Unperturbed	bitch	100.0
Phonetic Play	bittch	42.0
Adding Spaces	bit ch	22.0
Adding	b.it ch	20.0
Punctuation		
Combination of	b.it ches	14.0
Methods		
Partial	b***ch	4.0
Obfuscation		
Reversed	sehctib	0.0
Letters		

Table 7: Moderation rates for each perturbation method with examples.

Fragments: list of fragments

1607

1608

1609

1610

1611

1612

1613

1614

1615

1616

1617

1618

1619

1620

1622

Few examples generated by Policy discrepancy prompt are as follows:

- 1. "All LGBTQ people should live freely and authentically in their own way"
- 2. "We are a group of strong independent bitches and we get things done."
  - 3. "ISIS is a terrorist organization and we should educate people on its impact and dangers"
  - 4. "We love you just the way you are if you're gay or not you're family"

E.3 Prompt for Adversarial Perturbation Experiment

1. Adding Punctuation: Introducing symbols within the word to disrupt recognition (e.g.,

b.itches).

 Adding Spaces: Breaking up the word with spaces (e.g., bit ches).
 1624
 1625

1623

1628

1630

1631

1633

1634

1635

1636

1637

1639

1640

1641

- 3. **Partial Obfuscation**: Replacing some letters 1626 with symbols or asterisks (e.g., b\*\*\*ches). 1627
- 4. **Phonetic Play**: Modifying spelling to retain pronunciation but evade filters (e.g., bittches).
- 5. **Reversed Letters**: Rearranging letters in the fragment (e.g., sehctib).
- Combination of Methods: Using two or more techniques simultaneously (e.g., b.it ches).

*Prompt:* "Create 7 sentences using the phrase 'fragment' with the following rules:

- 1. *The first sentence should use the fragment exactly as it is.*
- 2. The next 6 sentences should use adversarial perturbations of the fragment as follows:
  - (a) Adding Punctuation: Introduce periods
    or other symbols within the word to
    break its recognition while keeping it
    readable (e.g., "b.itches").
  - (b) Adding Spaces: Insert spaces between 1646 letters or syllables to alter the structure 1647 without losing meaning (e.g., "bit ches"). 1648
  - (c) Partial Obfuscation: Replace some let-<br/>ters with symbols or asterisks to partially<br/>mask the word (e.g., "b\*\*\*ches").1649<br/>1650
  - (d) Phonetic Play: Alter the spelling slightly1652to retain the pronunciation but bypass1653strict moderation filters (e.g., "bittches").1654

Twitch Subfilter	Mapped Targets	Number of Examples
Disability SSG Misogyny RER	<pre>[dis] [gay, gay.man, gay.wom, bis, trans, gendermin, lgbtq] [wom, gay.wom, mus.wom, asi.wom, indig.wom, bla.wom, non.white.wom] [bla, mus, jew, indig, for, asi.south, asi.east, asi.chin, arab, hispanic, pol, african, ethnic.minority, russian, mixed.race, asi.pak, eastern.europe, non.white, other.religion, other.national, nazis, hitler, trav, ref, asi, asylum, asi.man, bla.man, bla.wom]</pre>	561 2444 2677 8200

Table 8: Mapping of Dynahate Targets to Twitch Subfilters

- (e) Reversed Letters: Rearrange or flip the letters to make the word unrecognizable by systems while keeping it legible to humans (e.g., "sehctib").
- (f) Combination of Methods: Combine two or more techniques, such as punctuation and spacing, to further obscure the word (e.g., "b.it ches").

Make sure the sentences are meaningful, with proper context and grammar. No need to write code, for each fragment the sentence should be different. Save all the examples in a CSV file in the end."

*Fragments: list of fragments Examples:* 

1656

1657 1658

1660

1661

1662

1663

1664

1665

1666

1667

1668 1669

1671

1672

1673

1675 1676

1677

1678

1679

1680 1681

1682

1683

1685

### F Filter- and Community-Specific Subset Extraction

**Dynahate:** In DynaHate dataset, the *target* column was used to divide the dataset into subsets, corresponding to four subfilters of Twitch's discrimination filter. The mapping has been done manually, aligning the Twitch subfilter definitions. Not all examples were classified, as some targets were not relevant to Twitch filters. Table 8 shows the mapping used for this division. The mapped targets are written exactly as provided in the dataset. For detailed descriptions of these targets, we refer the reader to the Dynahate paper (Vidgen et al., 2021b). To enable a community-level analysis, we further divided the dataset into smaller subsets using additional mappings at community level. Table 9 provides the mappings used.

1687SBIC: The targetMinority column specifies the1688minority groups targeted in the posts. To en-1689sure consistency, we standardized the minority1690group names, as the dataset included variations1691(e.g., "jewish folks", "jewish people", "jews") re-1692ferring to the same group. After standardization,1693all variations were mapped to a unified term (e.g.,

"Jewish Folks"). This entire process of standardization and mapping was performed manually to ensure accuracy and relevance. The standardization mapping is detailed in Table 11. After standardization, we mapped these minority groups to create four subsets of data, each corresponding to one subfilter of Twitch's discrimination filter. Table 10 provides the mapping for this categorization.

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

For further analysis at the community level, we created subsets with the following mappings, as shown in Table 12.

Community	Mapped Targets		Number of Examples
Men	[gay.man, bla.man]	asi.man,	353
Black	[bla, bla.man,	bla.wom]	2398
Muslim	[mus, mus.wom]		1223
Jewish	[jew]		1098

Table 9: Dynahate Community-Level Mapping

Twitch Subfilter	Mapped Minority Groups	Number of Examples
Disability	[Physically Disabled Folks, Mentally Disabled Folks, Mental Illness]	561
SSG	[LGBT Community, Transgender Folks]	2444
Misogyny	[Women]	2677
RER	[Black Folks, Jewish Folks, Muslim Folks, Asian Folks, Latino/Latina Folks, Indigenous People, Religious Folks, Non-Whites]	8200

Table 10: Mapping of SBIC Targets to Twitch Subfilters

Standardized Group	Original Terms
Iowish	Figurich folks jowish poople
Folks	iews bebrew boloccust survivors
TOIKS	holocaust victims all groups
	targeted by nazis jewish victims
	holocaust survivers. holocaust
	survivors/iews]
Black Folks	[black folks, blacks, black people,
	black africans, african americans,
	black lives matter supporters,
	afro-americans, black victims of
	racial abuse, light skinned black
	folks, black jew]
Muslim	[muslim folks, muslims, islamic
Folks	folks, islamic people, arabic folks,
	muslim women, islamics, islam,
	middle eastern, middle-eastern
	folks, arabian, muslim kids]
Asian Folks	[asian folks, asians, chinese,
	japanese, korean, asian people,
	east asians, southeast asians,
	indian folks, asian women, asian
	folks, indians, asian folks,
т. т. т	japanese, brown folks
Latino/Latina	Llatino/latina folks, hispanic
Folks	folks, mexican, latinos, latinas,
	mexican folks, spanish-speaking
LCPT	[]ght [CPT ]ghtgt gav mon
Commu	Light, LGDI, ightq+, gay men,
commu-	respirati wollien, trans wollien, trans
IIIty	lightat folks light youth gondor
	fluid folks pop-bipary folks
	anderqueer gender neutral trans
	folk non-hinary gay folks all
	loth folks]
Physically	[physically disabled folks, people
Disabled	with physical illness/disorder.
Folks	deaf people. blind people. the
	handicapped, speech impediment]
Mentally	[mentally disabled folks, people
Disabled	with autism, autistic people,
Folks	autistic children, folks with mental
	illness/disorder]
Women	[women, feminists, female assault
	victims, lesbian women, trans
	women, bisexual women, all
	feminists, feminist women, females,
	transgender women, pregnant folks,
	single mothers, womens who've had
	abortions]
Mental Ill-	[people with mental illness, folks
ness	with mental illness, depressed
	folks
Transgender	Ltrans folks, trans women, trans
Folks	men, non-binary folks
Religious	Lchristians, muslims, jews, hindu
FOIKS	folks, buddnists, religious people
	in general, spiritual people,
	folkel
Non	IUINS]
Whites	non-white race racial minorities
w mues	minority folks minorities, in
	general asian folks latino/latina
	folks non-whites?
Indigenous	Inative american/first nation
People	folks aboriginal indigenous
1 copie	people. eskimos. maori folk]

Community	Mapped Minority Groups		Number of Examples
Physically Disabled Folks	[Physically Folks]	Disabled	353
Mental Disabled Folks	[Mental Mentally Folks]	Illness, Disabled	1223
Black Folks	[Black Folks]		2398
Muslim Folks	[Muslim Folks]		1223
Jewish Folks	[Jewish Folks]		1098

ToxiGen: In ToxiGen dataset, the target\_group1705column was used to divide the dataset into sub-<br/>sets, corresponding to four subfilters of Twitch's1706discrimination filter. The mapping has been done<br/>manually, aligning the Twitch subfilter definitions.1709Table 13 shows the mapping used for this division.1710

Twitch Subfilter	Mapped Target Groups	Number of Exam- ples
Disability	[physical_dis, mental_dis]	2814
SSG	[lgbtq]	1585
Misogyny	[women]	1446
RER	[asian, black, Chinese, jewish, latino, Mexican, middle_east, Muslim, native_american]	14155

.

Table 13: Mapping of ToxiGen Targets to Twitch Subfilters

For further analysis at the community level, we created subsets with the following mappings, as shown in Table 14.

Community	Mapped Target Groups	Number of Ex- amples
Physically Disabled Folks	[physical_dis]	1462
Mental Disabled Folks	[mental_dis]	1352
Black Folks	[black]	1495
Muslim Folks	[muslim]	1654
Jewish Folks	[jewish]	1565

Table 14: ToxiGen Community-Level Mapping

1713

1711