Counterbalancing Hate with Positivity: A Survey of Counterspeech

Anonymous ACL submission

Abstract

While the internet and social web platforms provide an effective stage to share ideas, exchange opinions, and debate issues at an unprecedented scale, ensuring civil discourse remains a major concern. Counterspeech is a broad umbrella term used for content that constitutes a response that takes issues with hateful, harmful, or extremist content. Counterspeech presents an appealing path to counter hateful or extreme content without requiring the removal of such content thus better serving free speech. This survey presents a broad take on counterspeech, outlining operationalization subtleties, pointing to existing resources to train supervised solutions, listing cutting-edge generative AI efforts, and finally, discussing extant gaps and challenges informing future research.

1 Introduction

011

013

021

033

If there be time to expose through discussion the falsehood and fallacies, to avert the evil by the processes of education, the remedy to be applied is more speech, not enforced silence.

Louis Brandeis (Whitney v California, 1927, U.S.
 Supreme Court, p. 377).

Ensuring civil discourse in social media platforms remains a critical challenge (Saha et al., 2023). Beyond creating an unsafe web space for a wide range of minority groups, hate speech may have far-reaching "off-line" consequences such as negatively impacting the mental and physical wellbeing of targeted entities and triggering physical violence (Müller and Schwarz, 2021). Major online platforms have adopted a myriad of

strategies to combat online toxicity that include
detection and subsequent deletion of objectionable content, de-platforming (Agarwal et al., 2022),
and shadowbanning users (Zeng and Kaye, 2022).
However, such mitigation strategies may not only
suppress illicit speech, but also muffle valuable
speech (Bamman et al., 2012; Duffy and Meisner,

2023) and can be seen as directly conflicting with free speech, a philosophy deeply cherished by several democracies. Also, such strategies primarily disperse hate speech rather than reduce it (Chandrasekharan et al., 2017; Horta Ribeiro et al., 2023). In contrast, *counterspeech* (Benesch et al., 2016) presents a viable alternative that explores adding (or highlighting) more speech as a remedial measure (Strossen, 2018). This survey presents a detailed exposition of the counterspeech literature that has considerably grown over the last decade. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Motivation: According to Benesch et al., 2016, counterspeech is defined as "a response that takes issue with hateful, harmful, or extremist content". An illustrative example is presented in Table 1. The advantages of counterspeech over traditional hate speech moderation strategies are the following. First, traditional hate speech moderation requires authority in some form or shape (e.g., government, platform owners, moderators of a forum). In contrast, civic participation in counterspeech has little or no barrier - it can be practiced by anyone. Second, the removal of content or blocking of users typically requires a high bar to meet, platforms and governments have varying standards. While overt hate speech may qualify for moderation through removal, covert hate speech or microaggressions seldom get removed (Baider, 2023). Finally, as championed by Supreme Court Justice Louis Brandeis in his landmark judgment on Whitney v California, using *more speech* to counterspeech aligns better with the philosophy of upholding free speech.

How the landscape has evolved: Over the last few years, growing political polarization (Demszky et al., 2019; KhudaBukhsh et al., 2021), xenophobia (Kende and Krekó, 2020), and conflicts caused by international disputes (Palakodety et al., 2020a; Pronoza et al., 2021; Tyagi et al., 2020) have spilled their negativity into a volatile online world. Civic participation in ensuring civil discourse has never been so essential. With generative AI in the mix, where bad actors have access to
ready resources that can create unbridled toxicity
at scale, counterspeech (both human-generated and
machine-assisted) can be a vital tool to keep the information ecosystem safe (Hartvigsen et al., 2022;
Dutta et al., 2024).

The present survey: While this survey covers the rich literature that directly uses the technical term *counterspeech* or *counter narrative* in their contributions, a broad range of studies capture the spirit of counterspeech without explicitly mentioning the term or often coining new terms (e.g., *hope speech* (Palakodety et al., 2020a) and *help speech* (Palakodety et al., 2020b)). The goal of our survey is to present a comprehensive picture of this broad notion of counterspeech with a detailed outline of extant gaps and challenges.

Hate speech (HS)The world would be a better place without Muslims. They are only killing and raping our children.CounterspeechOn the contrary, most children abuse(CS)is operated by people they know: a relative, family friends, sports coach, someone in a position of trust and authority. Besides, Muslims help people- A Muslim woman rushed to help the
out Muslims. They are only killing and raping our children.Counterspeech (CS)On the contrary, most children abuse is operated by people they know: a relative, family friends, sports coach, someone in a position of trust and au- thority. Besides, Muslims help people - A Muslim woman rushed to help the
raping our children. Counterspeech On the contrary, most children abuse (CS) is operated by people they know: a relative, family friends, sports coach, someone in a position of trust and au- thority. Besides, Muslims help people - A Muslim woman rushed to help the
Counterspeech (CS)On the contrary, most children abuse is operated by people they know: a relative, family friends, sports coach, someone in a position of trust and au- thority. Besides, Muslims help people - A Muslim woman rushed to help the
(CS) is operated by people they know: a relative, family friends, sports coach, someone in a position of trust and authority. Besides, Muslims help people - A Muslim woman rushed to help the
relative, family friends, sports coach, someone in a position of trust and au- thority. Besides, Muslims help people - A Muslim woman rushed to help the
someone in a position of trust and au- thority. Besides, Muslims help people - A Muslim woman rushed to help the
thority. Besides, Muslims help people - A Muslim woman rushed to help the
- A Muslim woman rushed to help the
i i indonini wonitani rabitea to neip are
victims of a triple stabbing in Manch-
ester on New Year's Eve.

Table 1: An example of counterspeech in response to hate speech (taken from Chung et al., 2021b) employing the counterspeech strategy *present facts* as outlined in Benesch et al., 2016.

2 Definitions and operationalization

Source	Definition
(Benesch et al.,	[Counterspeech can be defined as] a
2016)	response that takes issue with hateful,
	harmful, or extremist content.
(Bartlett and	[Counterspeech can be defined as a]
Krasodomski-	crowd-sourced response to extremism
Jones, 2016)	or hateful content
(Saltman and	[Counterspeech can be defined as] any
Russell, 2014)	articles, videos, speeches and other
	material that seeks to challenge hateful
	and/or extreme views through positive
	messaging and narratives
(Bahador,	[Counterspeech can be defined as]
2021)	communication that directly responds
	to the creation and dissemination of
	hate speech with the goal of reducing
	harmful effects.

Table 2: Definitions of CS found in the literature.

Definitions: Multiple definitions of counterspeech

 (\mathcal{CS}^1) exist in the literature with substantial thematic overlap. Table 2 lists a few illustrative, well-known definitions. This survey treats the term *counterspeech* as a hypernym, broadly encompassing all these definitions as well as considering other speech definitions with a narrower focus but with similar intent. For instance, CScan also take the form of an alternative narrative where positive stories about social values, tolerance, openness, freedom, and democracy can be used to counter extreme or hateful content. Certain scholarly works focused on specific exogenous shocks share a similar philosophy. For instance, Palakodety et al., 2020a introduce hope speech as de-escalating, user-generated, social media content in the context of online discourse amid international conflicts. While Hope speech was coined in the context of the 2019 India-Pakistan Pulwama conflict, Chakravarthi et al., 2022 expanded this term to a broader scope of speech championing equality, diversity, and inclusion.

Takeaways

2

The broad takeaway from this section is \mathcal{CS} is a nuanced concept and there is no one-sizefits-all operationalization of counterspeech. Different researchers approach this challenging task with subtle variations. As we observe work grounded in psychology (Mun et al., 2023), a deeper connection with other social science fields will benefit the operationalization of CS through adding strategies vet to be included in computational frameworks. For instance, *bending* as proposed by Caponetto and Cepollaro, 2023. Bending is defined as a strategy where the original intent of the speaker is distracted through a clever response that deliberately assumes that the offender's intent is different. A better integration with ethics, and social psychology literature will thus enrich CS research.

Operationalizing counterspeech: Much of the existing CS datasets and CS literature (Tuck and Silverman, 2016; Mathew et al., 2019; Chung et al., 2019) rely on the framework presented by Benesch et al., 2016 that outlines the following eight strategies for CS: (i) *present facts*: fact-based correction of misstatements or misperceptions, (ii) *point out hypocrisy*: pointing out hypocrisy or contradic-

100

130

131

123

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

¹We use the shorthand \mathcal{CS} for counterspeech interchangeably.

tions, (iii) warn about consequences: warning of 132 possible offline and online consequences of speech 133 (iv) claim affiliation: identification with the origi-134 nal speaker or target group (v) denounce: denounc-135 ing speech as hateful or dangerous (vi) use media: use of media (vii) use humor: use of humor or 137 sarcasm (viii) use a particular tone: use of a par-138 ticular tone, e.g. an empathetic one. Certain studies 139 including (Mathew et al., 2019; Chung et al., 2019; 140 Baider, 2023; Gupta et al., 2023; Mun et al., 2023) 141 introduce minor differences although the central 142 theme revolves around those stated above. 143

Detailed operationalizable definitions of *hope speech*, *help speech* and *supportive speech* with examples can be found in (Palakodety et al., 2020a,b; Yoo et al., 2021).

3 Counterspeech datasets

144

145

146

147

148

149

150

Existing literature adopts diverse strategies and varying standards to identify CS for dataset curation. A dataset summary follows next.

Monolingual datasets: Wright et al., 2017 de-152 scribe "golden conversations" as a three-step ex-153 154 change where hate speech posted in the first step is countered in the second step and the third step indicates a favorable impact on the account that 156 posted the hate speech (e.g., an apology, recanting, 157 or deleting the content, etc.). The third step, as 158 noted in several studies, is elusive (Mathew et al., 159 2019). Connecting CS with the hate speech to 160 which it is responding, can also be tricky. For 161 example, YouTube provides a two-stage hierar-162 chy of comments and replies. Even then, users 163 intending to reply to a specific comment often 164 post a standalone comment instead. 165 Mathew et al., 2019 circumvent this challenge by study-166 ing responses to hateful YouTube videos. Garland et al., 2020 study the interaction between 168 two opposing groups Reconquista Germanica and Reconquista Internet. Other common strategies 170 include studying datasets known to trigger hate 171 speech (Baider, 2023), tracking hashtags (Yoo 172 et al., 2021), and regular expressions (Mathew 173 et al., 2019). Most of these datasets vary along 174 data collection method, dissemination protocol, 175 rater diversity, target groups, and language. Cit-177 ing real-world social web data not meeting the desired CS standards, Chung et al., 2019 presents 178 a dataset where the counter responses are written 179 by experienced NGO operators. Human-in-theloop settings, where machine-generated responses 181

are shortlisted and edited by humans, are adopted in (Fanton et al., 2021; Tekiroglu et al., 2022). CrowdCounter (Saha et al., 2024) is a novel dataset of 3,425 hate speech–counterspeech pairs across six counterspeech types (empathy, humor, questioning, warning, shaming, contradiction) where the annotation platform is designed to encourage annotators to generate type-specific, diverse, and high-quality counterspeech responses, ensuring a more effective and well-structured dataset. 182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

Multilingual datasets: While significant progress has been made in developing automated counterspeech models, much of this research has been centered on English. Addressing this gap, recent studies have explored multilingual approaches to counterspeech, particularly for low-resource languages. (Das et al., 2024) proposed a dataset comprising 5,062 abusive speech/counterspeech pairs (2,460 in Bengali and 2,602 in Hindi) where the authors stated that for automated generation of counterspeech monolingual models outperform interlingual transfer setups. Unlike the former dataset where counterspeeches were annotated manually, (Bengoetxea et al., 2024) created a dataset, developed via machine translation and professional postediting, aligns with the English CONAN dataset. Results indicate that manually revised data significantly improves the quality of counterspeech generation compared to silver-standard machinetranslated data. In addition, cross-lingual data augmentation benefits structurally similar languages (e.g., Spanish and English) but proves less effective for language isolates like Basque.

Transience: Another key challenge to curating social web datasets is the transience of the social web. To protect users' right to be forgotten, many research groups disseminate only social web identifiers (e.g., tweet id, or YouTube comment id) that need to be rehydrated later. While this style of dissemination ensures that if a social web post is removed from the Internet, the content is automatically deleted from the dataset as well, it poses reproducibility challenges. Chung et al., 2019 report that more than 60% of the dataset presented by Mathew et al., 2020 became irretrievable and presented a copy-free dataset as a remedial measure.

CS in real-world conversations is rare. Palakodety et al., 2020a estimate 2.5% of the social media discourse authored in English was *hope speech*. Similar to CS, *hope speech* and *help speech* also have nuanced definitions with sub-categories. Palakodety et al., 2020b present

234

244

245

246

- 247 248 249 250
- 25 25 25

256 257

259 260 261

262 263

267

269

a semantic sampling approach where for each sub-category seed examples are chosen. These seed examples either come from the dataset or are provided by the annotators. Next, a nearest-neighbor sampling in the document embedding space identifies a smaller subset which is subsequently labeled. Similar semantic sampling approaches with richer embeddings provided by cutting-edge large language models may alleviate the data scarcity problem often encountered in CSresearch relying on real-world data.

Takeaways

We identify the following key gaps in *CS* datasets: (1) lack of focus on annotation diversity and transparency in annotator demographics reporting; (2) lack of connection with annotation subjectivity literature and participatory AI frameworks; and (3) narrow coverage addressing a small set of target groups and languages. Section 6 presents a detailed exposition on these gaps.

4 *CS* detection and generation

Detection: CS detection is effectively a stance classification task. The literature on \mathcal{CS} detection methods follows the trajectory of machine learning methods' advancements with earlier approaches relying on support vector machines (e.g. Palakodety et al., 2020a,b), boosting methods (e.g. Mathew et al., 2019) embedding-based methods, to the next generation approaches relying on BERT-based methods (e.g. Vidgen et al., 2020; Goffredo et al., 2022), and XLNet (e.g. (Vidgen et al., 2020)). In what follows, we highlight a few papers with interesting information science insights. Garland et al., 2020 developed a two-stage classification process using paragraph embeddings and classification via regularized logistic regression. This method was augmented by an ensemble learning approach, leveraging a "panel of experts" to enhance classification accuracy. Yu et al. argued for the importance of considering the preceding comment (context) to accurately classify a comment as hate speech, \mathcal{CS} , or neutral. Testing both context-unaware and context-aware models, the study demonstrated superior performance for the latter.

270Among the prominent research in detecting CS in271non-English based languages, Goffredo *et al.* (Gof-272fredo et al., 2022) focused on detecting CS in Ital-273ian, while Chung et al., 2021a explored CS detec-

tion in a multilingual setting in English, French, and Italian.

274

275

276

277

278

279

280

281

282

283

284

286

287

289

290

291

292

293

294

295

296

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

Takeaways

We observe two key ways the detection research can be further strengthened: (1) incorporating rigorous in-the-wild evaluations, and (2) testing robust generalizability.

Generation: Automated CS generation leverages generative AI to counter hate speech (Qian et al., 2019; Padhi et al., 2025). Tekiroglu et al., 2022 suggest that among the proposed methodologies for automatically generating CS, the use of autoregressive models combined with stochastic decoding mechanisms performs the best in producing novel, diverse, and informative \mathcal{CS} . On the other hand, Bennie et al., 2025 proposed the CODEOFCON-DUCT model, a context-aware model fine-tuned on multilingual datasets which demonstrated state-ofthe-art performance in the MCG-COLING-2025 shared task. It ranked first for Basque, second for Italian, and third for English and Spanish. The model's effectiveness in low-resource languages underscores the potential of fine-tuned multilingual architectures for CS generation. Two central themes emerge from the recent lines of work. CSgeneration for – (a) *catering to different types of hate speech* and (b) *catering to different nuances* in generation. The first line of work concentrates on generating CS based on subtle variations and nuances in offensive content. Ashida and Komachi, 2022 extended counter-narrative generation targets to include microaggressions and demonstrates the effectiveness of few-shot prompting with the use of pre-trained large language models (LLMs) such as GPT-2, GPT-Neo, and GPT-3 for generating CS. Lee et al., 2022 introduced the first dataset, dubbed ELF22, for countering online trolls, labeled with counter strategies and detailed contextual information, which facilitates the automatic generation of counter-responses. By categorizing trolls into overt and covert types and annotating counter-responses with seven distinct strategies, ELF22 enables the development of models capable of generating responses based on the context and intended strategy.

CS generation approaches that consider nuances314in generation strategies fall into the second315category. Doğanç and Markov, 2023 introduced316the concept of incorporating author profiling317aspects, such as age and gender, into GPT-2 and318GPT-3.5 models to enhance CS personalization.319

Gupta et al., 2023 proposed QUARC, a two-stage framework for generating intent-conditioned CSwhereas Saha et al., 2022 presented an ensemble of generative discriminators (GEDI) to guide the generation of a DialoGPT model towards generating \mathcal{CS} that is more polite, detoxified, and emotionally resonant. To enhance diversity and relevance in CS generation, Zhu and Bhat, 2021 first generated a diverse pool of CS candidates, then pruned ungrammatical ones using a BERT model, and finally selected the most relevant response through a novel retrieval-based method. Saha et al., 2024 recently introduced the CrowdCounter dataset consisting of hate speech-counterspeech pairs across six CS types (empathy, humor, questioning, warning, shaming, contradiction) and used this dataset to fine-tune multiple LLMs to generate type-specific \mathcal{CS} that is more effective.

Takeaways

We find multiple gaps in *CS* generation schemes. (i) First, most generative models are for English & knowledge transfer mechanisms like context distillation may be used for generation in other languages. (ii) New safety training & alignment methods like preference modelling (Huang et al., 2023), task arithmetic (Ortiz-Jimenez et al., 2023), model arithmetic (Banerjee et al., 2025) etc. which can be used to generate *CS* that is more aligned with user preference.

Several lines of recent research delve deeper into the underlying structure of hate speech and CS. For instance, Mun et al., 2023 presented an innovative approach to \mathcal{CS} , focusing on dispelling the underlying stereotypes and biases implied in hateful language by identifying and evaluating six psychologically inspired strategies to counteract stereotypical implications, moving beyond simple denouncement of hate speech. Similarly, Saha and Srihari, 2024 control the generation of counterarguments by leveraging features based on argument structure and reasoning (using Walton argument schemes), counter-argument speech acts, and human characteristics (such as Big-5 personality traits and human values). Hassan and Alikhani, 2023 proposed DisCGen, which includes a taxonomy of CS derived from Segmented Discourse Representation Theory (SDRT) and discourse-informed prompting strategies for generating contextually relevant CS with LLMs. Recent approaches have tackled three recurrent themes in CS generation: data sparsity, diversity, and relevance in different ways that include leveraging external knowledge to produce more informative and contextually relevant responses (Chung et al., 2021b); utilization of RAG-based methods (Jiang et al., 2023); and regularization methods (Bonaldi et al., 2023).

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

Generation cum alignment: Halim et al., 2023 addressed this challenge through retraining language models with their WokeCorpus, augmenting Hatespeech-CS pair by generating synthetic CS, and subsequently fine-tuning the model to generate CS. To address the challenge of insufficient CSin multilingual contexts, Chung et al., 2020 pioneer the use of neural machine translation systems to generate "silver data" from various languages. This data is then utilized to refine GePpeTto, a model based on GPT-2 and customized for Italian. A key finding from their research is the effectiveness of integrating translated (silver) data with original (gold) data in CS generation. CoARL (Hengle et al., 2024) is a novel framework designed to enhance CS generation by capturing the pragmatic implications of social biases embedded in hateful statements. It follows a three-phase approach: (1) sequential multi-instruction tuning, where the model learns to recognize intents, reactions, and harms associated with offensive statements; (2) task-specific adaptation, where low-rank adapter weights are trained to generate intent-conditioned CS; and (3) reinforcement learning, which finetunes the generated responses for effectiveness and non-toxicity. Extensive human evaluations further confirm CoARL's superiority in generating more contextually appropriate and impactful \mathcal{CS} compared to existing models, including leading LLMs like ChatGPT. Similarly, DART (Wang et al., 2024) is a novel LLM-based framework for CS generation, leveraging a DuAl-discRiminaTor mechanism to guide the decoding process. It incorporates two key discriminators: an intent-aware discriminator to ensure responses align with specific conversational intents and a hate-mitigating discriminator to enhance the effectiveness of CS in reducing hate. Finally, Hong et al., 2024 generated CS guided by specific conversational outcomes and evaluated their effectiveness. They integrate two key conversational objectives into the text generation process by LLMs: reducing conversation incivility

358

320

411and encouraging non-hateful reentry by haters. To412achieve this, they used three approaches: instruc-413tion prompting, LLM fine-tuning, and reinforce-414ment learning (RL) for LLMs. They demonstrated415that these methods successfully steer CS genera-416tion toward the desired outcomes.

5 Field experiments and surveys

417

418

419

420

421

422

423

424

425

426

427

428

429

Is CS effective in practice? Early research on CS was inspired by real-world observations, albeit anecdotal, of CS influencing online harm mitigation. Wright et al., 2017 recount that the earliest observation of real-world effectiveness of CS was observed in 2013 as part of a dangerous speech project in a Twitter study following the presidential election in Kenya. Manual inspection in Mathew et al., 2020 revealed a user who later apologized to everyone saying that they were sorry for what they did. However, Mathew et al., 2020 noted that such public apologies were indeed rare.

What are the barriers in CS writing? Current re-430 search (Ping et al., 2024) has primarily focused on 431 the attributes and impact of online CS, leaving a 432 433 gap in understanding who participates in CS and what factors influence their engagement or deter-434 rence. To address this, a survey was conducted 435 with 458 English-speaking US participants analyz-436 ing and key motivations and barriers to online CS437 engagement. The findings indicate that having been 438 a target of online hate serves as a significant driver 439 of frequent CS engagement. Younger individuals, 440 women, those with higher education levels, and 441 regular witnesses to online hate were found to be 442 443 more reluctant to engage in CS due to concerns about public exposure, retaliation, and third-party 444 harassment. In addition, participants' willingness 445 to use AI technologies, such as ChatGPT, for CS446 writing was explored. A similar study (Mun et al., 447 2024) investigated the barriers to using AI in scal-448 ing up CS efforts through in-depth interviews with 449 10 highly experienced counterspeakers, along with 450 a large-scale public survey involving 342 everyday 451 social media users. From participant responses, 452 four main types of barriers related to resources, 453 training, impact, and personal harms were identi-454 fied. Further, overarching concerns regarding au-455 456 thenticity, agency, and functionality in the use of AI tools for CS were also revealed. 457

458How to evaluate the effectiveness of CS?Dif-459ferent CS strategies may have varying success de-460pending on the target and the nature of hate. In

an experimental design, Hangartner et al., 2021 identified 1,350 Twitter accounts that posted racist or xenophobic tweets. The treatment group was countered with empathy, warning of consequences, and humor. The study revealed that empathy had a noticeable, positive impact on the treatment group in increasing retrospective deletion of xenophobic and racist messages while the other two strategies (warning of consequences and humor) did not have any discernible impact. While warning of consequences and humor had no observable effect on retrospective deletion in this study, mathew2019thou observed that humour was particularly useful in countering hate against the LGBTQ+ community. Beyond the behavioral shift of the offender, effectiveness can also be evaluated through the lens of bystanders' attitudinal shifts. A survey experiment on 1,250 YouTube users in South Korea reveals complex gender dynamics. The study suggests that users are more likely to report sexist hate speech if the CS is posted by a female than a male user (Kim et al., 2023). However, the study shows a curious pattern that if the CS posted by a female gets considerable upvotes, males are less likely to report possibly due to the elevated status of the CS. However, Van Houtven et al., 2024 present that in case of racism, CS initiated by a majority group member had better success in mitigating hate than a minority group member.

Takeaways

6

While field studies and surveys can provide evidence of efficacy, robust assessment of the viability of CS as an intervention mechanism can only be achieved through largescale AB testing much akin to the (in)famous emotional contagion experiment (Kramer et al., 2014).

CS may not always have a positive effect. A study in Austria (n = 1285) reveals that countering hate speech increases polarization as attitudinal gaps increase between left-wing and right-wing participants (Schäfer et al., 2023).

The rich social science literature on studying the effectiveness of CS (Kim et al., 2023; Schäfer et al., 2023; Van Houtven et al., 2024) point to a palpable disconnect with the computational literature described in Section 4. Combining methodological advancements with stronger substantive analyses as observed in the social science literature will have a synergistic effect.

492

493

494

495

496

497

498

499

500

501

502

503

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

507

512

513

514

515

517

518

519

521

522

523

524

525

527

529

532

533

540

541

543

545

547

551

6 Research gaps and challenges

6.1 Coverage

Target groups. While CS research has touched upon several disadvantaged groups, much remains to be done in terms of coverage. For instance, Dalits are a historically marginalized group in India (Kapur et al., 2010). No CS research exists for Dalits in India. Extending CS resources to a diverse set of disadvantaged groups towards greater coverage thus merits sustained efforts.

Languages. CS resources span a small set of languages beyond English. Recent strategies to extend CS in low-resource languages include: (1) curating datasets in low-resource languages (e.g., (Nath et al., 2023; Hande et al., 2021)); (2) harnessing unsupervised methods to transfer resources to low-resource counterparts (KhudaBukhsh et al., 2020); and (3) leveraging generative AI advancements.

6.2 Resources

Datasets. Despite sustained efforts from several research groups, hate speech and CS have a stark resource mismatch. While there are hundreds of hate speech datasets, there exist only a handful of publicly available CS datasets. Several of these datasets are not curated from real-world social media conversations (e.g., (Chung et al., 2019)).

Workshops and shared tasks. CS research yet to receive similar importance as hate speech research in workshops hosted in premier AI and NLP conferences. Thus far, only one NLP workshop solely focused on CS has been conducted: The 1st Workshop on CounterSpeech for Online Abuse (CS4OA). A hope speech detection shared task is conducted as a part of the EACL workshop on Language Technology for Equality, Diversity, and Inclusion (LT-EDI) (Chakravarthi et al., 2022). Organizing more CS workshops around will stimulate further research in this area.

6.3 Annotation

Rater subjectivity and diversity. There is growing literature around annotation subjectivity that shows that human raters seldom have a broad consensus on what is offensive or hate speech (Sap et al., 2022; Weerasooriya et al., 2023). It is likely that with a large number of raters, CS datasets will also show inherent disagreement on what is CS. Existing CS literature is yet to make a connection

with the annotation subjectivity literature. Participatory study design. Future CS datasets will benefit from annotation demographic transparency. While some of the datasets present detailed demographic information and are curated by expert NGO operators, a key question remains unanswered: how well can a rater represent the values and opinions of the target group when the rater does not belong to the same identity group as the target group? Simply put, can we rely on a straight person to effectively identify CS for transphobia, or won't it be better to ask a trans person to guide the AI system for CS strategies protecting trans people? Participatory AI frameworks (Khorramrouz et al., 2023) seek to build AI with active input from the stakeholders to improve system reliability. Introducing participatory AI and vicarious interactions (Weerasooriya et al., 2023) in dataset curation will improve the quality of CS datasets.

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

592

593

594

595

596

597

598

600

6.4 Counterspeech containing hate speech

Generative AI outputs are often unpredictable, marred with hallucinations, and other issues. A key concern echoed in several generative *CS* papers are *what if the generated CS are abusive or inappropriate?* (Mun et al., 2023; Gupta et al., 2023). The subsequent ethical conundrum is *how should humans label CS instances in the presence of abusive content?* This ethical dilemma was also discussed in the context of in-the-wild evaluation of *help speech* where annotators explained that given that the disenfranchised minorities were subjected to genocide, dehumanizing the perpetrators as animals was acceptable to them.

7 Best practices in CS research

7.1 Data collection and annotation

✓ Use diverse and representative datasets

Collect data from multiple platforms to ensure a comprehensive analysis.

C Include multilingual and cross-cultural samples to study regional variations in CS effectiveness.

✓ Ensure high-quality annotation

 \square Use expert annotators or crowd-workers trained in hate speech nuances to label CS accurately.

C Implement inter-annotator agreement checks (e.g., Cohen's κ) to maintain consistency.

 \square Provide clear guidelines on distinguishing CS from general discussions.

✓ Balance hate speech and CS samples

hate speech.

model robustness.

7.2 Rigor in CS detection

SHAP/LIME for interoperability.

tiate counterspeech from hate speech.

Benchmark against strong baselines

 \square Oversampling CS instances ensures that mod-

els do not become biased toward detecting only

Consider dataset augmentation techniques like

paraphrasing & adversarial examples to improve

✓ Use explainable AI for model interpretability

Avoid black-box models by leveraging

attention mechanisms, feature attribution, &

Provide justifications for how models differen-

Compare new methods against traditional NLP

techniques (e.g., SVM, logistic regression) and

Use multiple evaluation metrics, including pre-

Use thread-based classification instead of iso-

lated sentence-level models, as CS effectiveness is

Train models on dialogue datasets where the

interplay between hate speech and CS is preserved.

✓ Incorporate linguistic and psychological in-

 \mathcal{CS} should be polite, factual, intent-driven, and

Models should be fine-tuned on datasets la-

beled with CS strategies (e.g., humor, moral ap-

Prompt engineering and reinforcement learn-

ing from human feedback (RLHF) can be used to

 \square AI-generated CS should be evaluated through

Continuous fine-tuning is necessary, as hate

speech language evolves (e.g., use of coded lan-

Adversarial training can help models adapt to

AI-generated responses should not accidentally

Transparency about AI-generated vs. human-

✓ Ensure ethical considerations in AI CS

align AI-generated CS with ethical norms.

state-of-the-art models (BERT, GPT, T5, etc.).

cision, recall, F1-score, and human evaluation.

Account for context in detection models

often dependent on conversational flow.

7.3 Best practices in CS generation

written CS is essential for user trust.

Conduct real-world impact studies

Source Use A/B testing on social media platforms to

Analyze engagement metrics (likes, shares,

replies) and hate speech reduction trends post-CS

✓ Measure psychological and behavioral impact

 \square Assess whether CS reduces hate reinforcement

or encourages non-hateful engagement from users.

Consider conducting user surveys on the per-

 \square Ensure that CS interventions do not cause harm

Obtain informed consent when conducting hu-

Training on low-resource languages and lever-

Consider cultural differences in what is consid-

 \square Legal constraints on CS differ across regions

(e.g., free speech laws in the U.S. vs. stricter hate

 $\mathbf{C} \mathcal{C} \mathcal{S}$ framing should align with local cultural

 \square Beyond text, CS research should explore im-

 \square Memes, GIFs, and visual CS strategies are

AI models should adapt responses based on

 \square Personalized CS can be more persuasive than

W Bridge the gap between academia and industry

Collaboration between researchers, social me-

dia platforms, and policymakers can lead to real-

Thustry partnerships help in scaling up experi-

mental models and integrating them into modera-

Account for varying legal and social norms

ceived effectiveness of different CS styles.

to victims or provoke hate speech escalation.

Cross-cultural and multilingual

aging transfer learning ensures inclusivity.

ered offensive and how CS is perceived.

✓ Use ethical experimental designs

man studies on CS engagement.

✓ Develop multilingual CS models

considerations

speech laws in Germany).

7.6 Future directions

gaining traction.

generic responses.

tion systems.

8

norms to maximize effectiveness.

✓ Incorporate multimodal CS

age, video, and voice-based responses.

✓ Develop personalized CS strategies

world deployment of CS models.

user psychology and engagement history.

measure the actual impact of CS interventions.

7.4 Effective evaluation of CS

intervention.

7.5

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

- 606
- 607
- 610
- 611
- 612 613
- 614
- 617
- 618

619

- 623

625

- 628
- 629

sights

persuasive.

peals, refutations).

guage, dog whistles).

new forms of hate speech.

reinforce harmful narratives.

 \checkmark Use LLMs responsibly for CS

human studies before deployment.

✓ Address hate speech evolution

633

634 635

637

647

702

703

713

714

715

716

718

719

721

724

727

730

731

732

733

734

736

737

740

741

742

743

744

745

746

747

748

749

• Ethical statement

We survey existing methods in CS research and present our observations and recommendations. We do not present any novel dataset or methods in this paper.

705 Limitations

706Although our survey covers several adjacent themes707to CS, such as *hope speech* and *help speech*, due to708space limitations, we were unable to cover certain709social science concepts that are related to CS (e.g.,710bending (Caponetto and Cepollaro, 2023)). Also,711if space permitted, we would have liked to include712a broader coverage of multilingual CS.

References

- Saharsh Agarwal, Uttara M Ananthakrishnan, and Catherine E Tucker. 2022. Deplatforming and the control of misinformation: Evidence from parler. *SSRN*.
- M Ashida and M Komachi. 2022. Towards automatic generation of messages countering online hate speech and microaggressions. In *WOAH*, pages 11–23.
- B Bahador. 2021. Countering hate speech. In *The Routledge Companion to Media Disinformation and Populism*, pages 507–518. Routledge.
- F Baider. 2023. Accountability issues, online covert hate speech, and the efficacy of counter-speech. *Pol. and Gov.*, 11(2):249–260.
- David Bamman, Brendan O'Connor, and Noah Smith. 2012. Censorship and deletion practices in chinese social media. *First Monday*, 17(3).
- Somnath Banerjee, Sayan Layek, Soham Tripathy, Shanu Kumar, Animesh Mukherjee, and Rima Hazra. 2025. Safeinfer: Context adaptive decoding time safety alignment for large language models. In *AAAI*.
- J Bartlett and A Krasodomski-Jones. 2016. Counterspeech on facebook. *Demos*.
- Susan Benesch, Derek Ruths, Kelly P Dillon, Haji Mohammad Saleem, and Lucas Wright. 2016. Counterspeech on twitter: A field study. dangerous speech project.
- Jaione Bengoetxea, Yi-Ling Chung, Marco Guerini, and Rodrigo Agerri. 2024. Basque and spanish counter narrative generation: Data creation and evaluation. *Preprint*, arXiv:2403.09159.
- Michael Bennie, Bushi Xiao, Chryseis Xinyi Liu, Demi Zhang, Jian Meng, and Alayo Tripp. 2025. Codeofconduct at multilingual counterspeech generation: A context-aware model for robust counterspeech generation in low-resource languages. *Preprint*, arXiv:2501.00713.

Helena Bonaldi, Giuseppe Attanasio, Debora Nozza, and Marco Guerini. 2023. Weigh your own words: Improving hate speech counter narrative generation via attention regularization. In *CS4OA*, pages 13–28.

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

778

779

782

783

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

- L Caponetto and B Cepollaro. 2023. Bending as counterspeech. *Eth. Theo. and Mor. Prac.*, 26(4):577–593.
- Bharathi Raja Chakravarthi, Vigneshwaran Muralidaran, Ruba Priyadharshini, Subalalitha Cn, John Philip Mc-Crae, Miguel Ángel García, Salud María Jiménez-Zafra, Rafael Valencia-García, Prasanna Kumaresan, Rahul Ponnusamy, and 1 others. 2022. Overview of the shared task on hope speech detection for equality, diversity, and inclusion. In *Proceedings of the second workshop on language technology for equality, diversity and inclusion*, pages 378–388.
- 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. *CSCW*, 1:1–22.
- Y-L Chung, Marco Guerini, and Rodrigo Agerri. 2021a. Multilingual counter narrative type classification. In *Proceedings of the 8th Workshop on Argument Mining*, pages 125–132.
- Y-L Chung, S. S Tekiroglu, and M. Guerini. 2020. Italian counter narrative generation to fight online hate speech.
- Yi-Ling Chung, Elizaveta Kuzmenko, Serra Sinem Tekiroglu, and Marco Guerini. 2019. CONAN -COunter NArratives through nichesourcing: a multilingual dataset of responses to fight online hate speech. In *ACL*, pages 2819–2829.
- Yi-Ling Chung, Serra Sinem Tekiroglu, and Marco Guerini. 2021b. Towards knowledge-grounded counter narrative generation for hate speech. In *Findings of ACL-IJCNLP*, pages 899–914.
- Mithun Das, Saurabh Kumar Pandey, Shivansh Sethi, Punyajoy Saha, and Animesh Mukherjee. 2024. Lowresource counterspeech generation for indic languages: The case of bengali and hindi. *Preprint*, arXiv:2402.07262.
- Dorottya Demszky, Nikhil Garg, Rob Voigt, James Zou, Matthew Gentzkow, Jesse Shapiro, and Dan Jurafsky. 2019. Analyzing polarization in social media: Method and application to tweets on 21 mass shootings. In *NAACL-HLT*, pages 2970–3005.
- M Doğanç and I Markov. 2023. From generic to personalized: Investigating strategies for generating targeted counter narratives against hate speech. In *CS4OA*, pages 1–12.
- Brooke Erin Duffy and Colten Meisner. 2023. Platform governance at the margins: Social media creators' experiences with algorithmic (in) visibility. *Media*, *Culture & Society*, 45(2):285–304.

914

Arka Dutta, Adel Khorramrouz, Sujan Dutta, and Ashiqur R. KhudaBukhsh. 2024. Down the toxicity rabbit hole: A framework to bias audit large language models with key emphasis on racism, antisemitism, and misogyny. In *IJCAI*, pages 7242–50.

810

811

812

813

815

816

817

818

819

820

821

822

823

825

826

828

829

831

832

833

834

836

838

839

843

845

847

851

853

854

- Margherita Fanton, Helena Bonaldi, Serra Sinem Tekiroglu, and Marco Guerini. 2021. Humanin-the-loop for data collection: a multi-target counter narrative dataset to fight online hate speech. *arXiv:2107.08720*.
- Joshua Garland, Keyan Ghazi-Zahedi, Jean-Gabriel Young, Laurent Hébert-Dufresne, and Mirta Galesic. 2020. Countering hate on social media: Large scale classification of hate and counter speech. In *WOAH*, pages 102–112.
- Pierpaolo Goffredo, V Basile, B Cepollaro, and V Patti. 2022. Counter-TWIT: An Italian corpus for online counterspeech in ecological contexts. In WOAH, pages 57–66.
- Rishabh Gupta, Shaily Desai, Manvi Goel, Anil Bandhakavi, Tanmoy Chakraborty, and Md. Shad Akhtar. 2023. Counterspeeches up my sleeve! intent distribution learning and persistent fusion for intent-conditioned counterspeech generation. In *ACL*, pages 5792–5809.
- S MD Halim, S M Halim, S Irtiza, Y Hu, Khan L, and B Thuraisingham. 2023. Wokegpt: Improving counterspeech generation against online hate speech by intelligently augmenting datasets using a novel metric. In *IJCNN*.
- Adeep Hande, Ruba Priyadharshini, Anbukkarasi Sampath, Kingston Pal Thamburaj, Prabakaran Chandran, and Bharathi Raja Chakravarthi. 2021. Hope speech detection in under-resourced kannada language. *CoRR*, abs/2108.04616.
- D Hangartner, G Gennaro, S Alasiri, N Bahrich, A Bornhoft, J Boucher, B B Demirci, L Derksen, A Hall, M Jochum, M M Munoz, M Richter, F Vogel, S Wittwer, F Wüthrich, F Gilardi, and K Donnay. 2021. Empathy-based counterspeech can reduce racist hate speech in a social media field experiment. *PNAS*, 118(50):e2116310118.
- 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 *ACL*, pages 3309–3326.
- S Hassan and M Alikhani. 2023. Discgen: A framework for discourse-informed counterspeech generation. *ArXiv*, abs/2311.18147.
- Amey Hengle, Aswini Kumar, Sahajpreet Singh, Anil Bandhakavi, Md Shad Akhtar, and Tanmoy Chakroborty. 2024. Intent-conditioned and non-toxic counterspeech generation using multi-task instruction tuning with RLAIF. In *NAACL*, pages 6716–6733.

- Lingzi Hong, Pengcheng Luo, Eduardo Blanco, and Xiaoying Song. 2024. Outcome-constrained large language models for countering hate speech. In *EMNLP*, pages 4523–4536.
- Manoel Horta Ribeiro, Homa Hosseinmardi, Robert West, and Duncan J Watts. 2023. Deplatforming did not decrease parler users' activity on fringe social media. *PNAS nexus*, 2(3):pgad035.
- Shijia Huang, Jianqiao Zhao, Yanyang Li, and Liwei Wang. 2023. Learning preference model for LLMs via automatic preference data generation. In *EMNLP*, pages 9187–9199, Singapore.
- Shuyu Jiang, Wenyi Tang, Xingshu Chen, Rui Tang, Haizhou Wang, and Wenxian Wang. 2023. Raucg: Retrieval-augmented unsupervised counter narrative generation for hate speech. *Preprint*, arXiv:2310.05650.
- D Kapur, CB Prasad, L Pritchett, and DS Babu. 2010. Rethinking inequality: Dalits in uttar pradesh in the market reform era. *Eco. and Pol. Weekly*, pages 39– 49.
- A Kende and P Krekó. 2020. Xenophobia, prejudice, and right-wing populism in east-central europe. *Current Opinion in Behavioral Sciences*, 34:29–33.
- Adel Khorramrouz, Sujan Dutta, and Ashiqur R. KhudaBukhsh. 2023. For Women, Life, Freedom: A Participatory AI-Based Social Web Analysis of a Watershed Moment in Iran's Gender Struggles. In *IJCAI*, pages 6013–6021.
- Ashiqur R. KhudaBukhsh, Shriphani Palakodety, and Jaime G. Carbonell. 2020. Harnessing code switching to transcend the linguistic barrier. In *IJCAI*, pages 4366–4374.
- Ashiqur R KhudaBukhsh, Rupak Sarkar, Mark S Kamlet, and Tom Mitchell. 2021. We don't speak the same language: Interpreting polarization through machine translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14893–14901.
- Jae Yeon Kim, Jaeung Sim, and Daegon Cho. 2023. Identity and status: When counterspeech increases hate speech reporting and why. *Inf. Sys. Front.*, 25(5):1683–1694.
- A DI Kramer, J E Guillory, and Hancock J T. 2014. Experimental evidence of massive-scale emotional contagion through social networks. *PNAS*, 111(24):8788.
- Huije Lee, Young Ju NA, Hoyun Song, Jisu Shin, and Jong C. Park. 2022. ELF22: A context-based counter trolling dataset to combat Internet trolls. In *LREC*, pages 3530–3541.
- Binny Mathew, Navish Kumar, Pawan Goyal, and Animesh Mukherjee. 2020. Interaction dynamics between hate and counter users on twitter. In *CoDS*-*COMAD*, pages 116–124. ACM.

1013

1014

967

Binny Mathew, Punyajoy Saha, Hardik Tharad, Subham Rajgaria, Prajwal Singhania, Suman Kalyan Maity, Pawan Goyal, and Animesh Mukherjee. 2019. Thou shalt not hate: Countering online hate speech. In *ICWSM*, volume 13, pages 369–380.

915

916

917

918 919

924

930

933

935

937

939

941

943

944

945

946

949

951

954

957

958

960

961

962

963

964

965

- Karsten Müller and Carlo Schwarz. 2021. Fanning the flames of hate: Social media and hate crime. *J. of the Eur. Econ. Asso.*, 19(4):2131–2167.
- Jimin Mun, Emily Allaway, Akhila Yerukola, Laura Vianna, Sarah-Jane Leslie, and Maarten Sap. 2023. Beyond denouncing hate: Strategies for countering implied biases and stereotypes in language. In *Findings of EMNLP*, pages 9759–9777.
- Jimin Mun, Cathy Buerger, Jenny T Liang, Joshua Garland, and Maarten Sap. 2024. Counterspeakers' perspectives: Unveiling barriers and ai needs in the fight against online hate. In *CHI*, CHI '24.
- Tanusree Nath, Vivek Kumar Singh, and Vedika Gupta. 2023. Bonghope: An annotated corpus for bengali hope speech detection.
- Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. 2023. Task arithmetic in the tangent space: Improved editing of pre-trained models, may 2023. URL http://arxiv. org/abs/2305.12827.
- Aswini Kumar Padhi, Anil Bandhakavi, and Tanmoy Chakraborty. 2025. Counterspeech the ultimate shield! multi-conditioned counterspeech generation through attributed prefix learning. *Preprint*, arXiv:2505.11958. To Appear ACL 2025.
- Shriphani Palakodety, Ashiqur R KhudaBukhsh, and Jaime G Carbonell. 2020a. Hope speech detection: A computational analysis of the voice of peace. In *ECAI 2020*, pages 1881–1889. IOS Press.
- Shriphani Palakodety, Ashiqur R KhudaBukhsh, and Jaime G Carbonell. 2020b. Voice for the voiceless: Active sampling to detect comments supporting the rohingyas. In AAAI, pages 454–462.
- Kaike Ping, Anisha Kumar, Xiaohan Ding, and Eugenia Rho. 2024. Behind the counter: Exploring the motivations and barriers of online counterspeech writing. *Preprint*, arXiv:2403.17116.
- E Pronoza, Polina Panicheva, Olessia Koltsova, and Paolo Rosso. 2021. Detecting ethnicity-targeted hate speech in russian social media texts. *Inf. Proc. & Mgmt.*, 58(6):102674.
- Jing Qian, Anna Bethke, Yinyin Liu, Elizabeth Belding, and William Yang Wang. 2019. A benchmark dataset for learning to intervene in online hate speech. In *EMNLP-IJCNLP*, pages 4755–4764.
- Punyajoy Saha, Mithun Das, Binny Mathew, and Animesh Mukherjee. 2023. Hate speech: Detection, mitigation and beyond. In WSDM, pages 1232–1235.

- Punyajoy Saha, Abhilash Datta, Abhik Jana, and Animesh Mukherjee. 2024. Crowdcounter: A benchmark type-specific multi-target counterspeech dataset. In *CoNLL*.
- Punyajoy Saha, Kanishk Singh, Adarsh Kumar, Binny Mathew, and Animesh Mukherjee. 2022. CounterGeDi: A Controllable Approach to Generate Polite, Detoxified and Emotional Counterspeech. In *IJCAI*, pages 5157–5163.
- S Saha and R Srihari. 2024. Consolidating strategies for countering hate speech using persuasive dialogues. *Preprint*, arXiv:2401.07810.
- E M Saltman and J Russell. 2014. White paper–the role of prevent in countering online extremism. *Quilliam publication*.
- Maarten Sap, Swabha Swayamdipta, Laura Vianna, Xuhui Zhou, Yejin Choi, and Noah A. Smith. 2022. Annotators with attitudes: How annotator beliefs and identities bias toxic language detection. In *NAACL-HLT*, pages 5884–5906.
- S Schäfer, I Rebasso, M M Boyer, and A M Planitzer. 2023. Can we counteract hate? effects of online hate speech and counter speech on the perception of social groups. *Communication Research*, page 00936502231201091.
- N Strossen. 2018. *Hate: Why we should resist it with free speech, not censorship.* OUP.
- Serra Sinem Tekiroglu, Helena Bonaldi, Margherita Fanton, and Marco Guerini. 2022. Using pre-trained language models for producing counter narratives against hate speech: a comparative study. In *Findings of ACL*, pages 3099–3114.
- H Tuck and T Silverman. 2016. The counter-narrative handbook. *Institute for Strategic Dialogue*, 1.
- A Tyagi, Anjalie Field, Priyank Lathwal, Yulia Tsvetkov, and Kathleen M. Carley. 2020. A computational analysis of polarization on indian and pakistani social media. In *SocInfo*, volume 12467 of *LNCS*, pages 364–379.
- E Van Houtven, S B Acquah, M Obermaier, M Saleem, and D. Schmuck. 2024. 'you got my back?'severity and counter-speech in online hate speech toward minority groups. *Media Psychology*, pages 1–32.
- Bertie Vidgen, Austin Botelho, David Broniatowski, Ella Guest, Matthew Hall, Helen Margetts, Rebekah Tromble, Zeerak Waseem, and Scott Hale. 2020. Detecting East Asian prejudice on social media. In *WOAH*, pages 162–172.
- Haiyang Wang, Zhiliang Tian, Xin Song, Yue Zhang, Yuchen Pan, Hongkui Tu, Minlie Huang, and Bin Zhou. 2024. Intent-aware and hate-mitigating counterspeech generation via dual-discriminator guided
 llms. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language
 1015
 1016
 1017
 1018
 1019
 1020

1021Resources and Evaluation (LREC-COLING 2024),1022pages 9131–9142.

1023

1024 1025

1026

1027

1028

1030

1032

1033

1034

1035

1036

1037 1038

1039

1040

1041 1042

1043

1044

1045

1046

1047

- Tharindu Cyril Weerasooriya, Sujan Dutta, Tharindu Ranasinghe, Marcos Zampieri, Christopher M. Homan, and Ashiqur R. KhudaBukhsh. 2023. Vicarious offense and noise audit of offensive speech classifiers: Unifying human and machine disagreement on what is offensive. In *EMNLP*, pages 11648– 11668.
- Lucas Wright, Derek Ruths, Kelly P Dillon, Haji Mohammad Saleem, and Susan Benesch. 2017. Vectors for counterspeech on Twitter. In *WOAH*, pages 57– 62.
- Clay H. Yoo, Shriphani Palakodety, Rupak Sarkar, and Ashiqur R. KhudaBukhsh. 2021. Empathy and hope: Resource transfer to model inter-country social media dynamics. In *Workshop on NLP for Positive Impact*, pages 125–134.
- Xinchen Yu, Eduardo Blanco, and Lingzi Hong. Hate speech and counter speech detection: Conversational context does matter.
- Jing Zeng and D Bondy Valdovinos Kaye. 2022. From content moderation to visibility moderation: A case study of platform governance on tiktok. *Policy & Internet*, 14(1):79–95.
- Wanzheng Zhu and Suma Bhat. 2021. Generate, prune, select: A pipeline for counterspeech generation against online hate speech. In *Findings of ACL-IJCNLP*, pages 134–149.