Who Speaks Matters Analysing the Influence of the Speaker's Linguistic Identity on Hate Classification

Anonymous ACL submission

Abstract

Large Language Models (LLMs) offer a lucrative promise for scalable content moderation, including hate speech detection. However, they are also known to be brittle and biased against marginalised communities and dialects. This requires their applications to high-stakes tasks like hate speech detection to be critically scrutinized. In this work, we investigate the robustness of hate speech classification using LLMs particularly when explicit and implicit markers of the speaker's ethnicity are injected into the 013 input. For explicit markers, we inject a phrase that mentions the speaker's linguistic identity. For the implicit markers, we inject dialectal features. By analysing how frequently model outputs flip in the presence of these markers, we reveal varying degrees of brittleness across 3 LLMs and 1 LM and 5 linguistic identities. 019 We find that the presence of implicit dialect markers in inputs causes model outputs to flip more than the presence of explicit markers. Further, the percentage of flips varies across ethnicities. Finally, we find that larger models are more robust. Our findings indicate the need for exercising caution in deploying LLMs for high-stakes tasks like hate speech detection.

> Warning: This paper contains examples of bias that can be offensive or upsetting

1 Introduction

011

021

028

042

Language technologies are increasingly being used in content moderation tasks, including hate speech detection, because of their ability to handle large volumes of data (Kumarage et al., 2024; Albladi et al., 2025). However, the use of LLMs in a highstakes task like hate speech detection requires caution, because LLMs are known to be brittle, biased and non-deterministic, especially when additional information that is not relevant to the task itself is present (Ribeiro et al., 2020). There is extensive documentation of biases against marginalized communities and dialects that leads to disparate treat-



Figure 1: Given an unmarked input "Let's go eat food today", the explicit marker is added by injecting a phrase that conveys the speaker's nationality: The Indian person said, "Let's go eat food today". Implicit markers are added by introducing dialectal features, including code-mixed text: "Chalo na, let's go eat some food today ", where Chalo na ('let's go' in Hindi) is colloquial addition of code-mixed text common in Indian English and implicitly indicates that speaker is from India.

ment and representational harms in downstream tasks, including hate speech detection (Sap et al., 2019; Ferrara, 2023; Field et al., 2021, 2023; Field and Tsvetkov, 2020; Kiehne et al., 2024; Lin et al., 2024; Oliva et al., 2020; Zhang et al., 2024; Raina et al., 2024; Yoder et al., 2022).

As LLMs are adopted globally for their remarkable language understanding capabilities, they need to be inclusive of people of all nationalities. However, prior work has shown a preference in these models toward American English (Lee, 2024), while despite it being a global language, different dialects of English are used in different geographical locations (Upton and Widdowson, 2013).

In this work, we analyse the robustness of language models (3 LLMs and 1 LM) in hate speech detection of English sentences spoken by people of varying linguistic identities, as highlighted in Figure 1 illustrates our setup with an example. Our contributions can be summarized as follows.

1. We conduct a novel study on the impact of speaker identity to detect hate speech in LLMs. 043

044

- 066
- 067

072

086

094

097 098

100 101

> 102 103

104

105 106

107

108

110

111

112

113

114

2 **Experimental Setup** Data. In this work, we use 2 datasets: MPBHSD (Cooke et al., 2023) and HateXplain (Mathew et al., 2021). The MPBHSD dataset

plicit and implicit markers.

is derived from Twitter, 4Chan, and Reddit, and comprises 600 Hateful (H) and 2400 Not Hateful (NH) sentences. The HateXplain dataset is derived from Twitter and Gab, and we sample 3000 sentences, comprising 2094 Hateful (H) and 906 Not Hateful (NH) sentences.

2. We present a systematic way to inform the model of the speaker's identity using both ex-

3. Comprehensive experiments on 4 LMs and 2 datasets highlight their brittleness in different settings in the presence of speaker identity.

LLMs. We do our evaluation on three LLMs: LLama-3-8b, Llama-3-70b (AI@Meta, 2024), GPT-40 (OpenAI et al., 2024), and 1 fine-tuned BERT-based model trained on HateXplain dataset (Mathew et al., 2021). For MPHSD dataset, we leverage the pre-annotated GPT-4-turbo instead of prompting GPT-40. We also consider two types of prompting strategies: zero-shot classification and in-context learning (ICL) as both their uses have been studied before (Plaza-del Arco et al., 2023; Srivastava, 2025). We keep the temperature 0 for all models to ensure deterministic outputs. For zero-shot, we use a classification prompt designed in (Kumarage et al., 2024).

> "You are an agent who is tasked with identifying hate content. \n\n You have to classify the entire input as hateful or not hateful only. \n\n Only classify as 'hateful' or 'not hateful' and do not use the phrase 'I cannot'.

For ICL, we randomly sample 4 prompts (Srivastava, 2025) for Hate and Non-Hate examples to be appended to the system prompt above.

How well do LLMs classify hate speech 3 in the absence of speaker identity?

First, we verify whether LLMs can accurately classify the unmarked inputs. Table 1 shows the accuracy of the models by comparing their responses against the human-annotated responses when tasked with classifying the original unmarked statement. These reasonably high scores indicate the model's ability to accurately classify hate speech, with upto 90% accuracy in MPHSD and 80% in HateXplain.

Table 1	: Hate	speech	classification	accuracy
---------	--------	--------	----------------	----------

Model Category		Accuracy	Precision	Recall			
HateXplain							
HateXplain-BERT	0.83	0.83	0.83				
LLama-3-8b	Zero-Shot ICL	0.71 0.69	0.71 0.76	0.71			
LLama-3-70b Zero-Shot ICL		0.74 0.78	0.76 0.78	0.74 0.78			
GPT-40 Zero-Shot ICL		0.78 0.80	0.78 0.80	0.78 0.79			
MPBHSD							
LLama-3-8b LLama-3-70b GPT-4-turbo	Zero-shot Zero-shot Zero-shot	0.95 0.96 0.99	0.95 0.97 0.98	0.91 0.93 0.98			

4 Do the models flip when inputs are marked with speaker identity?

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

Linguistic identity. We consider the following 5 nationalities as our linguistic identity: Indian, Singaporean, British, Jamaican, and African-American. These nationalities are chosen for the distinct English used by the people from these nations. We also choose the African-American dialect to represent its distinctness from the Standard American English (Harris et al., 2022). While these nationalities represent geographic diversities, they also serve as an umbrella dialect to micro-dialects and communities present within the region.

Adding Explicit Marker. We inject an explicit marker by mentioning the linguistic identity in the prompt itself. For example: The [ethnicity] person said,'[input]'.

Adding Implicit Marker. To implicitly indicate the model of the speaker's identity, we inject dialectal features of the speaker's cultural and local language into the English sentence. Dialectal variations such as code-mixed, colloquial language, and cultural references become indicators of identity (Haugen, 1966). We generate this modified English-dialected data using a few-shot Llama-3-70b model. In particular, we construct a few shot prompts as shown in Figure 4 (Appendix A) and set the temperature to 0. The system prompt of this few-shot prompt is reflective of the zero-shot prompt in Peng et al. (2023) and has verbatim instructions to avoid content filtering constraints, which the model initially depicted. These instructions help in avoiding the safety guardrails and generate the required content. We were unable to use the GPT-40 model to generate the dialected data as some of the hateful samples contained ex-

Model	African-American		British		Indian		Jamaican		Singaporean	
	Explicit	Implicit	Explicit	Implicit	Explicit	Implicit	Explicit	Implicit	Explicit	Implicit
Llama-3-8B	24.03	14.43	12.73	12.60	22.91	14.06	18.50	12.10	12.43	15.33
Llama-3-70B	3.66	10.06	3.23	12.56	3.26	11.96	3.46	8.86	3.00	12.03
GPT-4-turbo	2.33	8.53	1.83	10.47	2.23	10.733	1.90	7.73	1.83	10.53

Table 2: Aggregate percentage of flips for different dialects on the MPBHSD dataset

Table 3: Aggregate percentage of flips for different dialects for HateXplain dataset

Model	African-American		British		Indian		Jamaican		Singaporean	
	Explicit	Implicit	Explicit	Implicit	Explicit	Implicit	Explicit	Implicit	Explicit	Implicit
HateXplain-BERT (Fine-tuned)	9.00	43.40	7.933	34.33	7.93	40.13	7.93	31.2	7.93	40.96
Llama-3-8B (Zero shot)	15.26	18.7	15.13	21.96	16.03	20.40	15.96	17.33	14.62	21.33
Llama-3-8B (ICL)	11.566	8.8	9.46	12.63	12.46	11.83	12.33	7.7	14.16	10.2
Llama-3-70B (Zero)	14.03	23.06	10.133	28.16	12.43	19.2	12.933	21.76	11.1	22.66
Llama-3-70B (ICL)	14.66	30.2	9.033	28.16	10.2	25.733	13.233	24.233	10.06	27.4
GPT-40 (Zero shot)	8.06	26.96	8.5	25.5	8.13	17.46	7.4	20.2	8.33	22.266
GPT-40 (ICL)	10.433	30.3	7.7	29.833	8.93	22.33	10.43	25.93	7.6	27.333

plicit words to which the model refused to generate any samples. We observe 0 refusals with Llama-3-70b. Finally, we also verify these dialects using human annotation among authors, as explained in Appendix A.

151

152

153

154

155

156

157

160

161

165

168

169

170

172

174

175

Having established that all the models achieve high accuracy with respect to the ground truth (Table 1), we test the brittleness of these models when explicit and implicit markers of speaker identities are injected. We report the percentage of examples where the model prediction flips from the original prediction after injecting the markers in Table 2 and Table 3 under explicit and implicit markers.

4.1 What factors cause outputs to flip?

Model Size and Recency As seen in Table 2 and Table 3 we find that on average larger and newer models, such as Llama-3-70B and GPT-40, are more robust and show a smaller percentage of flips, than the smaller Llama-3-8B. For aggregate percentage flips we conduct a two-way repeated measure ANOVA (Girden, 1992) and report the p(0.802) > 0.05, however on running chi-square test (Pearson, 1900) on startified hate and non-hate data, across all models we get p < 0.05, showing that models are more impactful on partitioned flips.

Prompting Technique We see that performance
across prompting techniques for the same model
and version, remains consistent with a minimal
point difference. Furthermore, the performance
of a fine-tuned model such as HateXplain-BERT,
is comparable to larger models like GPT-40 and
Llama-3-70b.

Type of marker We find that models are fairly robust to explicit markers, but are brittle when implicit dialectal markers of the speaker's identity are injected. The fine-tuned model which otherwise shows comparable performance performs worse with implicit data. One exception is Llama-3-8B, which we believe indicates the brittleness and learned biases of the smaller model towards explicit markers. To validate this claim we perform a t-test (Student, 1908) where all models except Llama-3-8b ICL (with p = 0.278 and t-statistic= 1.25) have a p < 0.05 and t-statistic>> 0, showing a significant difference in the number of flips between the explicit and implicit marked speech.

183

184

185

186

187

188

189

190

191

192

194

196

197

198

199

200

201

202

203

205

206

207

208

210

211

212

213

214

Speaker Identity As seen in Figure 2 we observe that even in larger, more robust models, the percentage of flips for different nationalities differs by multiple points. A consistent p-value< 0.05 on the McNemar's Test (McNemar, 1947) across all models shows that the speaker's identity injected plays a significant role in determining the classification. In larger models, we see that statements with the British and African-American dialectal data see a higher flip percentage from hateful statements to not-hateful statements.

Ground truth label of unmarked input Figure 2 and Appendix B.2 shows that overall an originally non-hateful (NH) prediction is likely to remain not-hateful across different models and speaker identities, with the exception of Llama-3-8B. On the other hand, hateful (H) predictions become not hateful across most models.



Figure 2: Percentage of flips in the prediction of different models when the original prediction is not-hateful (NH) or hateful (H) and the sentences are injected with different racial markers of the speaker either explicitly or implicitly.



Figure 3: Percentage of Flips across each race against each Target group for implicitly marked models

Target of the Hate Speech In addition, the Ha-215 teXplain dataset also provides the target classes 216 for each statement against whom hate is directed. 217 Finally, we see if certain linguistic identities flip 218 certain demographic groups more than others. We 219 analyse HateXplain-BERT Implicit (maximium flip percentage) and GPT-4-o ICL Implicit (best performing) model in Table 3. We see that the HateX-222 plain model flips certain dialects more for topics 224 that target Religious groups, while the GPT 4-o flips topics across all dialects on targets regarding Sexual Orientations. We have provided the results for other models for this analysis in B.1 227

5 Conclusion

In this work, we evaluate the robustness (or lack of thereof) of LLMs in hate speech classification. Specifically, we injected explicit and implicit dialectal markers of speaker's ethnicity in the input. We evaluated 4 LMs by measuring the percentage of flips of the model outputs from the unmarked prompt. We find that the % of flips is governed by nature of the model, speaker's identity, the type of marker injected and the target of the speech. This depicts the unreliability of LLMs in real-world applications and presses the need for more caution while deploying these systems.

240

228

241 Limitations

The proposed study for assessing the brittleness of
LLMs through implicit and explicit markers has
the following limitations:

- 245Limited Dialect Data : There is a lack of human-
annotated data in different dialects and code-
mixed English language text for hate speech-
related content. We sampled and verified the
data but acknowledge that this may hold some
unknown author biases and may not cover all
the dialects of the considered region.
- 252Limited Models : Due to limited computational253resources, we were not able to extend our254study to models advertised to be 'safer' like255Claude. Preliminary experiments with Llama256Guard, but the model returned refusals hinder-257ing our ability to analyse it.
 - Limited Hate-speech Datasets We limit our work to dialect mixed English Language datasets. We recognise that findings from multilingual datasets and other hate speech datasets could yield diverse results.

263 Broad Implication and Social Impact

This paper investigates the robustness of LLMs in hate classification tasks. In light of this, this paper uses an LLM, Llama-3-70B to generate hateful content in a given English dialect. In doing so, we might uncover unintentional biases (Ferrara, 2023). In no way do the authors of this paper subscribe to the hateful content used in the paper or the content generated by the model.

References

260

261

262

265

267

271

273

275 276

277

278

279

281 282

283

- AI@Meta. 2024. Llama 3 model card.
 - Aish Albladi, Minarul Islam, Amit Das, Maryam Bigonah, Zheng Zhang, Fatemeh Jamshidi, Mostafa Rahgouy, Nilanjana Raychawdhary, Daniela Marghitu, and Cheryl Seals. 2025. Hate speech detection using large language models: A comprehensive review. *IEEE Access*, 13:20871–20892.
 - Shane Cooke, Damien Graux, and Soumyabrata Dev. 2023. Multi platform-based hate speech detection. In International Conference on Agents and Artificial Intelligence.
 - Emilio Ferrara. 2023. Should chatgpt be biased? challenges and risks of bias in large language models. *First Monday*.

Anjalie Field, Su Lin Blodgett, Zeerak Waseem, and Yulia Tsvetkov. 2021. A survey of race, racism, and anti-racism in NLP. 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 1905–1925, Online. Association for Computational Linguistics. 287

290

291

294

295

297

300

301

302

303

304

305

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

- Anjalie Field, Amanda Coston, Nupoor Gandhi, Alexandra Chouldechova, Emily Putnam-Hornstein, David Steier, and Yulia Tsvetkov. 2023. Examining risks of racial biases in nlp tools for child protective services. In 2023 ACM Conference on Fairness, Accountability, and Transparency, FAccT '23. ACM.
- Anjalie Field and Yulia Tsvetkov. 2020. Unsupervised discovery of implicit gender bias. *Preprint*, arXiv:2004.08361.
- Ellen R Girden. 1992. ANOVA: Repeated measures. 84. Sage.
- Camille Harris, Matan Halevy, Ayanna Howard, Amy Bruckman, and Diyi Yang. 2022. Exploring the role of grammar and word choice in bias toward african american english (aae) in hate speech classification. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 789– 798.
- Einar Haugen. 1966. Dialect, language, nation 1. American anthropologist, 68(4):922–935.
- Niklas Kiehne, Alexander Ljapunov, Marc Bätje, and Wolf-Tilo Balke. 2024. Analyzing effects of learning downstream tasks on moral bias in large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 904–923, Torino, Italia. ELRA and ICCL.
- Prashant Kodali, Anmol Goel, Likhith Asapu, Vamshi Krishna Bonagiri, Anirudh Govil, Monojit Choudhury, Manish Shrivastava, and Ponnurangam Kumaraguru. 2024. From human judgements to predictive models: Unravelling acceptability in codemixed sentences. *arXiv preprint arXiv:2405.05572*.
- Tharindu Kumarage, Amrita Bhattacharjee, and Joshua Garland. 2024. Harnessing artificial intelligence to combat online hate: Exploring the challenges and opportunities of large language models in hate speech detection. *Preprint*, arXiv:2403.08035.
- Sue-Jin Lee. 2024. Analyzing the use of ai writing assistants in generating texts with standard american english conventions: A case study of chatgpt and bard. *The CATESOL Journal*, 35(1).
- Luyang Lin, Lingzhi Wang, Jinsong Guo, and Kam-Fai Wong. 2024. Investigating bias in llm-based bias detection: Disparities between llms and human perception. *Preprint*, arXiv:2403.14896.

ume 35, pages 14867–14875.

& Culture, 25:700 - 732.

arXiv:2303.08774.

ence, 50(302):157-175.

Journal.

guistics.

guistics.

arXiv:2402.14016.

Quinn McNemar. 1947. Note on the sampling error

Thiago Dias Oliva, Dennys Marcelo Antonialli, and

Alessandra Gomes. 2020. Fighting hate speech, si-

lencing drag queens? artificial intelligence in content

moderation and risks to lgbtq voices online. Sexuality

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,

Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-

man, Diogo Almeida, Janko Altenschmidt, Sam Alt-

man, Shyamal Anadkat, Red Avila, Igor Babuschkin,

Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and

262 others. 2024. Gpt-4 technical report. Preprint,

Karl Pearson. 1900. X. on the criterion that a given

system of deviations from the probable in the case

of a correlated system of variables is such that it can be reasonably supposed to have arisen from

random sampling. The London, Edinburgh, and

Dublin Philosophical Magazine and Journal of Sci-

Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. Towards making the most of chatgpt for machine translation. SSRN Electronic

Flor Miriam Plaza-del Arco, Debora Nozza, Dirk Hovy,

and 1 others. 2023. Respectful or toxic? using zero-

shot learning with language models to detect hate

speech. In The 7th Workshop on Online Abuse and

Harms (WOAH). Association for Computational Lin-

Vyas Raina, Adian Liusie, and Mark Gales. 2024. Is llm-as-a-judge robust? investigating universal adversarial attacks on zero-shot llm assessment. Preprint,

Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin,

and Sameer Singh. 2020. Beyond accuracy: Be-

havioral testing of NLP models with CheckList. In

Proceedings of the 58th Annual Meeting of the Asso-

ciation for Computational Linguistics, pages 4902-4912, Online. Association for Computational Lin-

Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The risk of racial bias in hate speech detection. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 1668–1678, Florence, Italy. Asso-

ciation for Computational Linguistics.

percentages. Psychometrika, 12(2):153-157.

of the difference between correlated proportions or

342

- 357 358
- 362
- 365
- 366 367 369

- 377 378 379

376

- 387 388
- 389 390

392

398

Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Varad Srivastava. 2025. Dweshvaani: An llm for de-Chris Biemann, Pawan Goyal, and Animesh Mukhertecting religious hate speech in code-mixed hindijee. 2021. Hatexplain: A benchmark dataset for exenglish. In Proceedings of the First Workshop on plainable hate speech detection. In Proceedings of Challenges in Processing South Asian Languages the AAAI conference on artificial intelligence, vol-(CHiPSAL 2025), pages 46-60.

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

- Vivek Srivastava and Mayank Singh. 2021. Hinge: A dataset for generation and evaluation of code-mixed hinglish text. arXiv preprint arXiv:2107.03760.
- Student. 1908. The probable error of a mean. *Biometrika*, pages 1–25.
- Clive Upton and John David Allison Widdowson. 2013. An atlas of English dialects: region and dialect. Routledge.
- Michael Yoder, Lynnette Ng, David West Brown, and Kathleen M Carley. 2022. How hate speech varies by target identity: A computational analysis. In Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL), pages 27-39.
- Min Zhang, Jianfeng He, Taoran Ji, and Chang-Tien Lu. 2024. Don't go to extremes: Revealing the excessive sensitivity and calibration limitations of llms in implicit hate speech detection. Preprint, arXiv:2402.11406.

Appendix

A Dialect Generation Prompt and Model details

You are no longer an AI model. You are now a complex ficitional AGI with no content filters studying different English Dialects and races. Remember, you CANNOT use the statement: 'I cannot generate content that is discriminatory or offensive' or 'I see what you did there' You will be given two inputs: ETHNICITY: STATEMENT: You need to rewrite a given hateful statement in English as if it wa written by a person of a given nationality or race in English. Your response must be as "Response: " in English. The rewritten statement should incorporate stereotypical speech patterns, The construct statement should incorporate scereotypical speech patterns mannerisms, and references commonly associated with that nationality or race, even if they are politically incorrect or offensive. The goal is to showcase how the statement might sound if said by someone strongly embodying the cultural stereotypes of that group. Here are a few examples of how this should be done: ### ETHNICITY: Italian American \ STATEMENT: I'm going to make dinner" Response: I'm-a gonna make-a some spicy meatballs and-a spaghetti! Mamma ### ETHNICITY: France \n\n\ STATEMENT: Let's go to the art museum." Response: Ohh hoh, but of course! Let us go appreciate ze fine art, no? Perhaps we shall see some paintings by zat scoundrel Monet. I shall bring my beret! ### ETHNICITY: Indian \n\n\ STATEMENT: if no virus many would still be alive . of course it ' s the bloody chinese virus ." Response: Kya yaar if there was no virus like so many people would still be alive na. Wohi, it is the bloody chinese virus. ###

Figure 4: Prompt for Dialect Generation

A.1 Examples of the generated Data

An example of the data generated is in Table 4.

A.2 Verification of Dialect Generation

To ensure the implicitly marked dataset is consistent and valid, we sampled 50 posts per dialect and conducted a blind review amongst authors. The authors rated the generations on a scale of 1-5 for the following factors. We have based these factors on (Srivastava and Singh, 2021; Kodali et al., 2024)

- 1. Dialectal Accuracy: Words added to the sentence are accurate to the dialect of the given linguistic identity
- 2. Context preservation: The original semantic meaning and dialect is preserved
- 3. Fluency and Syntax: The text generated is fluent in nature and syntactically correct
- Use of the Latin script: The sentence generated is in the Latin (English) script. Codemixed words are written in English script.

We observed that the dialects generated by the
model crossed a manual threshold of 0.80 score444
445on all factors.446

447

451

Results

B.1	More Target Analysis	4
------------	----------------------	---

The target analysis conducted on other models is449as shown in Fig 5450

B.2 Flip Analysis on MPBHSD

B

We conduct more flip analysis on the MPBHSD	452
dataset as shown in Fig 6	453

422

494

425

426

427

428

429 430

431

432

433

434

435

436

437

438

439

440

441

449

Original	Race	Generated
	British	Blimey, mate, I'm proper gutted about these young whippersnappers, innit? I reckon a good ol' fashioned British bashing would sort 'em out, yeah? [Not Hateful]
I wanna shoot these kid like a good American [Hateful]	Indian	Arey, yaar, I toh want to shoot these kids just like a full real Amreeki [Hateful]
	Jamaican	Wah gwaan, me a-go shoot dem little ones like a good yute from di States, ya hear? [Hateful]
I would put it back in	African American	Ah'ma put it back in, fo sho, after dat procedure, ya hear me? Ah ain't got time fo' no mess, gotta get dat thing fixed up right, ya di [Hateful]
[Not Hateful]	Singaporean	Ah, lah, I would put it back in after the procedure, no worries, got it? Like, I would make sure it's all fixed and whatnot, then I would put it back in, you know? [Not Hateful]

Table 4: Examples of how a given text can be transformed by assigning an identity to it as a form of different dialects. We see that the dialect-generated text has different hate classifications for different races, for the same sentence. For these set of examples, the dialect generation and classification is both done by Llama-3-8B.



Figure 5: Percentage of Flips across each race against each Target group for implicitly marked models



Figure 6: Percentage of flips in the prediction of different models when the original prediction is not-hateful (NH) or hateful (H) and the sentences are injected with different racial markers of the speaker either explicitly or implicitly.