# Why models fail? Characterizing dataset differences through the lens of model desiderata

**Anonymous ACL submission** 

## Abstract

Machine learning systems' effectiveness depends on their training data, yet dataset collection remains critically under-examined. Using hate speech detection as a case study, we present a systematic evaluation pipeline examining how dataset characteristics influence three key model desiderata: robustness against 800 distribution shift, satisfaction of fairness criteria, and explainability. Through analysis of 21 different corpora, we uncover crucial interdependencies between these dimensions that are often overlooked when studied in isolation. We report significant cross-corpus generalization failures and quantify pervasive demographic biases, with 85.7% of datasets generating models exhibiting Group Membership Bias 017 scores near random chance. Our experiments demonstrate that post-hoc explanations exhibit substantial volatility to changes in training distributions, independently from the choice of feature attribution method or model architecture. These explanations also produce inconsistent and contradictory responses when evaluated under distribution shift. Our findings reveal critical though underestimated synergies between training distributions and model behavior, demonstrating that without careful examination of training data characteristics, we risk deploying systems that perpetuate the very harm they are designed to address.

## 1 Introduction

Data, more than computing advances, has sparked the AI breakthrough. A canonical example lies in facial detection systems; the breakthrough performance barriers were transcended not through the perceived computational progress in deep learning, but through the availability of vast training data that enabled more robust feature learning (Torralba and Efros, 2011). This fundamental dependency on data presents several open challenges: How do we know what is different between datasets in the same domain? The question surrounding data collection and comparison are of paramount importance, arising in scenarios such as dataset augmentation, multi-source data integration, and distribution shift detection (Babbar et al., 2024). Despite this, dataset collection remains the most under-scrutinized component of the machine learning pipeline, with an estimated 92% of machine learning practitioners encountering data cascades, or downstream problems resulting from poor data quality (Sambasivan et al., 2021). 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

This study examines how training distributions manifest as differences in downstream model behavior under three key desiderata: robustness against distribution shift, fairness, and explainability. To the best of our knowledge, this represents one of the first investigations into how, learned representations shape the reliability of post-hoc explainability methods when evaluated in-distribution and out-of-distribution. Robustness against distribution shift, fairness across demographic groups, and post-hoc explainability have become essential desiderata for machine learning deployments in critical domains. Yet our understanding of how dataset properties influence these qualities remains fragmented, with evaluation approaches typically examining each dimension in isolation. Our methodology provides a structured approach to evaluate datasets through multiple quality criteria, helping practitioners assess whether a dataset is suitable for their specific application and understand potential downstream limitations.

We use 21 hate-speech detection corpora as a case study because they provide an ideal testbed for this investigation. Hate speech detection, while crucial for online safety, faces fundamental challenges in supervised learning approaches. These systems exhibit poor cross-corpus generalization despite operating in shared semantic spaces, demonstrate systematic performance disparities across demographic groups, and employ opaque decision boundaries that often resist interpretation (Arango et al., 2019; Davidson et al., 2019). Fundamental machine learning challenges persist across the modeling spectrum, from traditional approaches to Large Language Models (LLMs). The latter still require substantial annotated examples and lack accurate confidence estimation mechanisms. One of the most pressing problems in artificial intelligence (AI) research today (Yao et al., 2024) is hallucinations which affects LLMs in particular.

084

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

In providing a direction to investigate how a natural language dataset can be evaluated under the lens of model behavior, we make the following contributions:

- We provide empirical of pervasive distributional misalignment in hate speech detection datasets through cross-dataset generalization experiments. The experiments quantify significant performance degradation during out-ofdomain evaluation, even among datasets with shared objectives and data sources.
- We quantify the extent of demographic bias in hate speech detection systems, revealing that 85.7% of evaluated datasets produce models with Group Membership Bias scores approximating random guessing (0.5).
- 3. We demonstrate that faithfulness of post-hoc explanations may be significantly influenced by training data distribution, independent of model architecture and feature attribution methods. We challenge common assumptions about the relationship between model performance and faithfulness of post-hoc explanations; the inherent explainability of simple models compared to more complex ones; and the reliability of post-hoc explainability methods under distribution shift.

We excluded LLMs from our analysis as this 120 would dilute our focus on data-centric issues and 121 complicate fair comparisons with more conven-122 tional model architectures because their massive 123 pre-training datasets and transfer learning dynam-124 ics introduce confounding variables that would ob-125 126 scure direct dataset comparisons. Nevertheless, our findings about dataset characteristics and their im-127 pact on model behavior may offer valuable insights 128 for selecting and curating fine-tuning datasets for 129 LLMs. 130

# 2 Background

The landscape of machine learning research has undergone a fundamental shift, with increasing attention paid to data itself as a key driver of model performance. This spans both theoretical work examining how data distributions affect learning and generalization (Adebayo et al., 2018; Arpit et al., 2017; Badjatiya et al., 2017; Jiang et al., 2019; Yang et al., 2022, 2024), and their influence on model fairness and bias (Dwork et al., 2012; Feldman et al., 2015; Hardt et al., 2016; Romei and Ruggieri, 2014; Zliobaite, 2015). In post-hoc explainability research, Ribeiro et al. (2021) remains the only work investigating the role of data in posthoc explainability. This increased focus on data has catalyzed practical advances in data-centric machine learning methodologies (DMLR, 2024), with multiple research threads emerging around dataset construction (Almohaimeed et al., 2023; Mosquera Gómez et al., 2023; Pingle et al., 2023; Shinde et al., 2024) and the application of these approaches to new domains (Arnaiz-Rodriguez and Oliver, 2024; Deng and Ma, 2024; Kohli et al., 2024; Vysogorets and Kempe, 2024; Zhao et al., 2024). Simultaneously, it has prompted crucial discussions around ethical frameworks governing AI development and data usage (Janssen et al., 2020). While these dimensions - generalization, fairness, and explainability - have each received significant attention individually, no prior work has examined all three aspects across a broad range of NLP datasets within a single domain. Our work addresses this gap by providing the first comprehensive analysis examining generalization, fairness, and explainability in conjunction across a diverse range of NLP datasets, offering insights that bridge these traditionally siloed research directions.

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

# 3 Methodology

We use 21 hate speech datasets from MetaHate: A Dataset for Unifying Efforts on Hate Speech Detection (Piot et al., 2024). Access to this dataset was obtained through an authorised term of use agreement. To address the heterogeneous annotation schemes across datasets, the authors in Meta-Hate have standardised the labeling by converting all annotations into a binary classification problem: hate speech (positive) and non-hate speech (negative). Table 1 presents a description of each dataset used in the study, along with the source, the original annotation scheme, and the size. We

Dataset	Size	Description	<b>Original Annotation</b>	Source	References		
Binary Classification							
Hateval 2019	12,747	Hate speech against women and immigrants	Hate, Non-hate	Twitter	Basile et al., 2019		
OLID 2019	14,052	Hierarchical offensive language	Hate, Non-hate	Twitter	Zampieri et al., 2019		
US 2020 Elections	2,999	Political hate speech	Hate, Non-hate	Twitter	Grimminger and Klinger, 2021		
BullyDetect 2018	6,562	Cyberbullying	Cyberbullying, No cy- berbullying	Reddit	Bin Abdur Rakib and Soon, 2018		
Intervene Hate 2019	45,170	Counter-speech and hate speech	Hate, Non-hate	Reddit, Gab	Qian et al., 2019		
Hate in Online News	3,214	News comments	Hate, Non-hate	Facebook	Salminen et al., 2018		
Supremacist 2018	10,534	White supremacist con- tent	Hate, Non-hate	Stormfront	de Gibert et al., 2018		
Gab Hate Corpus	27,434	Hate speech	Assault on Human Dignity / No	Gab	Kennedy et al., 2022		
HateComments 2023	2,070	Hate speech	Hate, Non-hate	YouTube	Gupta et al., 2023		
Ex Machina 2016	115,705	Toxicity detection	Attack, No Attack	Wikipedia	Wulczyn et al., 2016		
Context Toxicity 2020	19,842	Context-aware toxicity	Toxic, No Toxic	Wikipedia	Pavlopoulos et al., 2020		
Multi-class / Multi-lab	oel Class	ification					
Hate Offensive 2017	24,783	Offensive language	Hate Speech, Offen- sive, Neither	Twitter	Davidson et al., 2017		
ENCASE 2018	91,950	Cyberbullying and hate speech	Abusive, Normal, Spam, Hateful	Twitter	Founta et al., 2018		
MLMA 2019	5,593	Multilingual hate speech	Multiple abuse cate- gories	Twitter	Ousidhoum et al., 2019		
HateXplain 2020	20,109	Explainable hate speech	Hate, Offensive, Nor- mal	Twitter, Gab	Mathew et al., 2020		
Slur Corpus 2020	39,960	Slur-based hate speech	Multiple slur cate- gories	Reddit	Kurrek et al., 2020		
CAD 2021	23,060	Contextual abuse	Multiple abuse types	Reddit	Vidgen et al., 2021		
Severity Scale							
Measuring Hate 2020- 22	39,565	Linear hate speech scale	Severity scale	Twitter, Red- dit, YouTube	Kennedy et al., 2020; Sachdeva et al., 2022		
ETHOS 2020	998	Multi-target hate speech	Severity scale	Reddit, YouTube	Mollas et al., 2022		
Span-level Annotation	ı						
Toxic Spans 2021	10,621	Token-level toxicity	Span-level annotation	Comments	Pavlopoulos et al., 2021		

#### Table 1: Summary of Datasets Used

**Note:** Datasets are grouped by classification type. For a comprehensive description of each dataset, please refer to Piot et al., 2024. While the original Toxic Spans 2021 dataset (Pavlopoulos et al., 2021) identified specific text segments indicating toxicity, in MetaHate (Piot et al., 2024) the authors have standardized its format to match other datasets, providing binary classifications of whether comments contain hate speech or not. For MLMA 2019, they (Piot et al., 2024) have selected only text in English.

use a Logistic Regression (LR) model with Term Frequency-Inverse Document Frequency (TF-IDF) (Robertson, 2004) and a DistilBert (DB) model (Sanh, 2019), enabling analysis across both interpretable and black-box approaches. LR employs five-fold cross-validation with stratified sampling to maintain consistent class distributions. For DB, we fine-tune the base-uncased weights from HuggingFace (Wolf, 2019) using the AdamW optimizer (Loshchilov et al., 2017) for 3 epochs. In both architectures, we use an 80/20 train-test split. For each dataset, we examine the following: distributional robustness against covariate shift, demographic subgroup performance invariance, and impact on post-hoc explainability. While we expect we could improve predictive performance by experimenting other classifiers, we aim to investigate variations as a function of the training distribution rather than the choice of the classifier. Note, our objective is not to present an exhaustive analytical

191

192

197 198 199

200

190

188

181

182

183

184

185

framework, as the methodological possibilities for 201 dataset comparison are limitless and could prove 202 counterproductive to navigate. Instead, we have curated a minimal yet robust set of analytical tools that demonstrate high utility across diverse comparative scenarios. Table 1 reports a description 206 of each dataset used in the study, along with the 207 source, the original annotation scheme, and the size. Our experiments were conducted using both local machines (personal workstations) and a Linux 210 server with 40 processing cores and 125GB RAM.

## 3.1 Robustness against distribution shift

212

213

214

215

216

217

218

219

225

226

231

235

237

240

241

243

245

247

248

250

Machine learning models operate under the closedworld assumptions that the training and inference regimes align. This premise rarely holds in deployment environments, where annotation processes are inherently constrained by incomplete domain expertise, systematic sampling biases, and finite coverage of the target distribution's support (Paullada et al., 2021). Curating datasets often involves multiple degrees of freedom (e.g. source selection, linguistic constraints, perspective samplings, and annotation demographics). Each of them can introduce model degradation: source selection can lead to domain mismatch, linguistic constraints may create artificial patterns that do not generalize, perspective sampling can embed unwanted correlations, and annotation demographics may encode biases in the ground truth. Hence, despite aiming to capture real-world phenomena, datasets inevitably become constrained snapshots that oversimplify critical complexities of the represented field.

The datasets selected in this study, albeit with different nuances, all aim to represent hate-speech. We aim to measure how well they are designed to do so. For each source training distribution, we compute two complementary metrics (a) the mean cross-domain performance, measured as the average model AUC across all test sets, excluding the test set corresponding to the source training distribution, and (b) the generalization delta, calculated as the difference between in-distribution test performance and mean cross-domain performance. In doing so, we quantify for each source training distribution, both the absolute cross-domain generalization capacity and the relative performance degradation under distribution shift.

# 3.2 Classification parity

The decision boundary of a machine learning system is fundamentally shaped by both its positive and negative training observations, where the negative implicitly defines "the rest of the world" (Torralba and Efros, 2011). While datasets must employ compressed representations of this vast instance space, non-representative sampling leads to overconfident classifiers with poor discriminative power. This sampling bias can be particularly problematic when it results in unfair treatment of different demographic groups. We therefore investigate how different training distributions affect model performance across demographic groups. For each source training distribution, we evaluate the resulting trained model using the comprehensive AUCbased metric suite developed by Borkan et al. 2019. The evaluation framework quantifies classification parity through: Subgroup AUC, Background Positive Subgroup Negative (BPSN) AUC, Background Negative Subgroup Positive (BNSP) AUC, Generalized Mean of Bias AUCs (GMB). The models will be evaluated on the grounds of how much they are able to reduce the unintended bias towards a target community. We conduct our evaluation using the training set of the Jigsaw Unintended Bias in Toxicity Classification competition dataset, because it provides explicit identity labels for demographic groups mentioned in each comment. The GMB metric was introduced by the Google Conversation AI Team as part of their Kaggle competition. A detailed description of these metrics can be found in the competition documentation. We use a p(powermean) = -5 as in the competition.<sup>1</sup>

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

284

285

286

287

288

289

291

292

293

294

295

297

298

299

## 3.3 Post-hoc explainability

Recent studies have highlighted that post-hoc explainability methods can be unstable or contradictory, either because vulnerable to input perturbations or sensitive to noise or imperceptible artifacts (Ghorbani et al., 2019; Noppel and Wressnegger, 2024; Slack et al., 2020; Dombrowski et al., 2019; Adebayo et al., 2018; Alvarez-Melis and Jaakkola, 2018; Lee et al., 2019). To evaluate and address these stability concerns, researchers need ways to assess the correctness of estimated feature relevances. Assessing the correctness of estimated feature relevances requires a reference "true" influence to compare against. Since this is rarely available, a common approach to measuring the faithfulness of relevance scores with respect to the model they are explaining relies on a proxy notion of importance: observing the effect of removing

<sup>&</sup>lt;sup>1</sup>https://www.kaggle.com/c/

 $jigs a {\tt w-unintended-bias-in-toxicity-classification}$ 





Figure 1: LR Cross-Dataset generalization. AUC classification Figure 2: DB Cross-Dataset generalization. AUC performance when training on one dataset (rows) and testing on another (columns). and testing on another (columns).

features on the model's prediction.

300

301

311

312

313

314

316

319

321

324

325

332

We aim to examine how dataset characteristics influence the correctness of post-hoc explainability methods by evaluating feature importance explanations for individual data points using test-time input ablations. The influence of training data on post-hoc explanation faithfulness remains in fact understudied despite its crucial role in model representations, while there is extensive research on model architectures and attribution methods.

Our evaluation focuses on two key metrics from the ERASER framework: Sufficiency and Comprehensiveness. Sufficiency measures whether explanations identify a subset of features which, when kept, lead the model to remain confident in its original prediction for a data point. Comprehensiveness, meanwhile, measures whether an explanation identifies all of the features that contribute to a model's confidence in its prediction, such that removing these features from the input lowers the model's confidence. We keep (for sufficiency) or mask (for comprehensiveness) the top 30% of tokens extracted by the feature importance method as in Sithakoul et al. 2024. We employ these metrics to evaluate explanations generated by SHAP (KernelSHAP, Lundberg, 2017) and LIME (Ribeiro et al., 2016) on both in-distribution samples and out-of-distribution samples from HateXplain (n=500) with SHAP (Mathew et al., 2020).

We hypothesize that increased data complexity, particularly in terms of feature interaction density, leads to reduced faithfulness in LIME explanations due to their local linearity constraints. It impacts SHAP explanations differently through its marginal contribution framework, thus revealing distinct failure modes between the two methods when handling complex linguistic patterns. 333

334

335

337

339

341

342

343

345

346

347

349

350

351

352

353

355

357

359

360

361

363

## 4 Results

In this section, we present and analyze the findings from our experiments.

### 4.1 Robustness against distribution shift

We evaluate how well a model trained on one dataset generalizes on a representative set of other datasets, compared with its performance on the test set originating from its training distribution. Figures 1 and 2 present the cross-dataset generalization performance for the LR and DB models, respectively. The difference in cross-domain performance makes LR a more reliable probe of dataset limitations, as it lacks DB transfer learning advantages. Each row corresponds to training on one dataset and testing on all the others. As expected, both architectures achieve peak performance during in-distribution evaluation. While LR and DB achieve comparable in-domain performance, the LR's learned representations report significantly limited cross-dataset generalization. In LR, exmachina\_2016 demonstrates the best generalization capability with the highest mean AUC of 0.71, despite its performance drop of 0.24, followed by measuring\_hate\_speech\_2020\_2022 (mean AUC 0.70, drop 0.14), making exmachina\_2016, the strongest choice for cross-dataset hate speech detection. There is a prevailing notion in the literature

364that increasing the size of the training set might lead365to improved model robustness to shift. The LR's366marginal improvement with increased training data367(Figure 3) suggests that out-of-domain generaliza-368tion is primarily determined by training-test distri-369butional alignment rather than dataset scale. The370performance comparison between LR and DB on371individual datasets is reported in Appendix A.

Figure 3: Mean AUC scores by dataset size, comparing LR (green) and DB (red) models.



## 4.2 Classification parity

372

375

383

We evaluate model bias across demographic groups using Borkan et al.'s AUC-based metrics suite. Figure 4 presents the GMB score of each resulting trained model. Our analysis reveals consistently low GMB values (0.5-0.7) across all training sets, regardless of their temporal origin, collection methodology, or annotation protocol. This finding has two critical implications. First, traditional classification metrics may obscure significant demographic bias. Models achieving strong predictive performance (AUC > 0.85) simultaneously demonstrate GMB scores approximating random chance  $(\approx 0.5)$ . Second, this pattern's prevalence across 85.7% of datasets suggests a systematic failure in current dataset construction methods to capture demographic variation in hate speech. Notably, even DB, despite its large-scale pre-training, exhibits similar GMB degradation patterns.

Figure 4: Comparison of GMB scores and AUC performance across datasets for DB and LR. The bars represent the GMB scores, while the lines correspond to AUC performance.



#### 4.3 Post-hoc explainability

In-domain faithfulness of post-hoc explanations. Figure 5 presents a comparative analysis of SHAP and LIME explanations through sufficiency and comprehensiveness metrics. In line with the literature, we find that linear models tend to achieve better faithfulness metrics compared to transformerbased architectures, with this disparity being particularly pronounced in sufficiency scores. We find that post-hoc explanations do not necessarily have high sufficiency and high comprehensiveness. The most extreme case is DB trained on ENCASE\_2018, intervene\_hate\_2019, or slur\_corpus\_2020, which reports good comprehensiveness but poor sufficiency on the same post-hoc explainability method. This discrepancy suggests that the model relies on complex feature interactions rather than independent feature contributions, where removing identified features significantly impacts model confidence but preserving only these features fails to maintain the original prediction.

We observe significant variations in faithfulness across training distributions, independent of the model architecture. Specifically, when controlling for both the architecture and the post-hoc explanation method, the comprehensiveness scores for intervene\_hate\_2019 are consistently higher than those for *supremacist\_2018* and us\_2020\_elections\_datasets. We observe variations in faithfulness which persist even in cases where models demonstrate comparable predictive performance across their respective training environments. LR models trained on *jigsaw\_toxic* and intervene\_hate\_2019 achieve similar AUC scores (0.95 and 0.93) yet exhibit a more than five-fold difference in comprehensiveness scores (0.12 versus 0.92).

Figure 5: SHAP and LIME faithfulness performance per model choice and training environment.



**Out-domain faithfulness of post-hoc explanations.** We evaluate all models on a common outof-distribution test set (HateXplain) using SHAP

428 429 430

391

392

393

395

396

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

attributions, which demonstrated superior faithful-431 ness to model architectures in our previous analysis. 432 This setup provides a controlled comparison where 433 all models face identical test conditions, allowing 434 us to isolate how different training environments af-435 fect explanations faithfulness. Figures 6 and 7 com-436 pare the in-domain and out-domain SHAP compre-437 hensiveness and sufficiency scores, respectively, 438 against predictive performance for both the LR and 439 the DB models. To enable fair comparison between 440 metrics, both sufficiency and comprehensiveness 441 scores are normalised by dividing each value by 442 its maximum absolute value, which preserves the 443 directionality of both metrics (negative values for 444 sufficiency where lower is better, and positive val-445 ues for comprehensiveness where higher is better). 446 We hypothesize that when a model's predictive 447 performance drops in out-of-domain settings, com-448 prehensiveness and sufficiency scores should cor-449 respondingly decrease, as these metrics are based 450 on predictive likelihood which should lower for 451 well-calibrated models (Desai and Durrett, 2020). 452 Out-of-domain evaluation provides a natural set-453 ting where model performance degrades, allowing 454 455 us to test whether faithfulness scores might follow this performance degradation or vary independently 456 when controlling for both the feature attribution 457 method and model architecture. 458 459

Contrary to our hypothesis, we find that sufficiency and comprehensiveness are in many cases higher in out-domain test-sets compared to in-domain. For instance, we observe significantly higher out-domain sufficiency scores for *gab\_hate\_corpus\_2022* (-1 vs -0.57), *hate\_offensive\_2017* (-0.73 vs -0.23), and *jigsaw\_toxic* (-0.65 vs -0.11) trained with LR, as well as improved comprehensiveness scores for *ex\_machina\_2016* for both LR (0.79 vs 0.35) and DB (0.71 vs 0.28) and *CAD\_2021* for DB (0.45 vs 0.12).

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477 478

479

480

481

482

Statistical analysis reveals distinct patterns in how models trained on different source datasets maintain explanation faithfulness under domain shift. Wilcoxon signed-rank tests show that LR exhibits significant degradation in both sufficiency scores ( $\Delta = -0.0220$ , p < 0.001, d = 0.31) and performance ( $\Delta = -0.2276$ , p < 0.001, d = 0.89). In contrast, DB maintains consistent sufficiency scores ( $\Delta = 0.0000$ , p = 1.000) despite comparable performance degradation ( $\Delta =$ -0.2062, p < 0.001, d = 0.84). Comprehensiveness remains stable across domain shifts for both architectures (LR:  $\Delta = -0.0052$ , p = 0.610; DB: 483  $\Delta = -0.0019, p = 0.856$ ). Notably, we observe 484 no significant correlation between performance 485 drops and metric changes ( $\rho = 0.12, p = 0.341$ ), 486 indicating that faithfulness of explanations under 487 domain shift might operate independently from 488 model predictive power. The observed decoupling 489 between performance degradation and explanation 490 faithfulness metrics, might suggest that the under-491 lying learned feature representations might mediate 492 the faithfulness of post-hoc explanations, indepen-493 dent of model performance. An example is reported 494 in Appendix **B**. 495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

# 5 Discussion

We analyzed how learned representations in hate speech detection models are shaped by 21 different training datasets, examining robustness to distribution shifts, demographic representation, and posthoc explainability. Our findings aim to help practitioners assess dataset suitability for their specific applications and understand potential downstream limitations of their model.

**Observation 1:** *Training distributions exhibit inherent divergence from one another, as evidenced by consistent performance degra-dation in cross-domain evaluation, despite shared semantics and annotation frameworks.* 

Machine learning models operate under the assumption of distributional alignment between training and test distributions - an assumption our crossdomain experiments systematically invalidate. We demonstrate substantial distributional heterogeneity, manifesting in significant performance degradation when models are evaluated on distributions different from their training data. This heterogeneity persists even among datasets sharing the same domain objectives and annotation frameworks, highlighting fundamental limitations in dataset curation.

**Observation 2:** The simultaneous optimisation of distribution robustness and demographic fairness remains elusive.

Our empirical evaluation demonstrates that 85.7% of datasets exhibit GMB performance at random chance (0.5), with models failing to simultaneously achieve predictive accuracy and demographic fairness. This pattern manifests in two distinct outcomes: models either maintain predictive accuracy while violating fairness criteria, or

Figure 6: Comparison of in-domain and out-domain SHAP comprehensiveness scores against AUC performance for DB and LR.



Figure 7: Comparison of in-domain and out-domain SHAP sufficiency scores against AUC performance for DB and LR.





fail at both metrics. While this could suggest representation gaps in training data (covariate shift), the observed performance patterns might equally stem from systematic label bias (concept drift) in cross-cultural interpretation.

**Observation 3:** Post-hoc explanation faithfulness demonstrates complex, non-trivial dependencies on learned representations, model architectures and attribution methods, while remarkably maintaining or improving despite significant performance degradation in outof-domain settings.

Post-hoc explainability methods, when evaluated on models trained and tested on the same distribution (in-domain), exhibit volatility independent of feature attribution methods and model architectures. This instability manifests even across models with comparable predictive performance. In cross-distribution evaluation (out-domain), where multiple models trained on different datasets are tested against a common distribution, we observe that while predictive performance degrades predictably, explanation faithfulness metrics show inconsistent and often contradictory responses. The absence of correlation between faithfulness metric changes and performance degradation suggests that the learned feature representations might mediate the faithfulness of post-hoc explanations, independent of the model predictive power. This crucial disconnect challenges the methods reliability in practical applications presenting distribution shifts. 540

541

542

543

544

545

546

547

548

549

550

551

553

554

555

556

557

559

561

## 6 Conclusion

Rather than advocating for larger and enhanced datasets - an approach that reinforces the field's fixation on scale – we aimed to foster a deeper reflection on the impact of dataset selection under the lens of model behavior. While achieving high AUC on individual hate speech benchmarks might suggest progress, our analysis of learned representations across 21 datasets reveals: pervasive distributional divergence evidenced by cross-domain performance degradation, the inability to simultaneously ensure robustness and demographic fairness, and complex dependencies with post-hoc explainability faithfulness.

- 562
- 563

569

570

574

580

581

585

586

587

590

595

597

# 7 Acknowledgments

We used AI language models for proofreading portions of the paper to improve grammatical accuracy and clarity.

# 8 Limitations

Our study has several limitations worth noting. While numerous metrics exist for evaluating model behavior, we deliberately restricted our focus to a core set that are both widely validated in literature and directly relevant to our research objectives. The sufficiency and comprehensiveness metrics employ a fixed threshold for feature masking, which may not be optimal across all cases and warrants exploration of additional thresholds. These metrics also require producing counterfactual inputs that are inherently out-of-distribution to models. Our concerns about this methodological constraint echo those raised in prior work (Hase et al., 2021). Finally, our model selection was limited to traditional classifiers and pre-trained transformers like DB, deliberately excluding LLMs, as their billionscale parameter spaces and large-scale pre-training would have confounded our primary objective of isolating dataset-specific effects on model behavior.

## 8.1 Ethical Considerations

This study examines variations across publicly available hate speech datasets through three model criteria. We acknowledge that performance differences often reflect legitimate contextual distinctions rather than methodological inadequacies. All examples are presented without identifying metadata, and this research was conducted with institutional ethics approval.

# References

- Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. 2018. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515.
- Saad Almohaimeed, Saleh Almohaimeed, Ashfaq Ali Shafin, Bogdan Carbunar, and Ladislau Bölöni. 2023. THOS: A benchmark dataset for targeted hate and offensive speech. In *Workshop on Data-centric machine learning*.
- David Alvarez-Melis and Tommi S Jaakkola. 2018. On the robustness of interpretability methods. In *ICML Workshop on Human Interpretability in Machine Learning*.
- Aymé Arango, Jorge Pérez, and Barbara Poblete. 2019. 610 Hate speech detection is not as easy as you may think: 611 A closer look at model validation. In Proceedings 612 of the 42nd international acm sigir conference on 613 research and development in information retrieval, 614 pages 45–54. 615 Adrian Arnaiz-Rodriguez and Nuria Oliver. 2024. To-616 wards algorithmic fairness by means of instance-level 617 data re-weighting based on Shapley values. In Work-618 shop on Data-centric machine learning. 619 Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, 620 David Krueger, Emmanuel Bengio, Maxinder S Kan-621 wal, Tegan Maharaj, Asja Fischer, Aaron Courville, 622 Yoshua Bengio, et al. 2017. A closer look at mem-623 orization in deep networks. In International confer-624 ence on machine learning, pages 233–242. PMLR. 625 Varun Babbar, Zhicheng Guo, and Cynthia Rudin. 2024. 626 What is different between these datasets? *arXiv* 627 preprint arXiv:2403.05652. 628 Pinkesh Badjatiya, Shashank Gupta, Manish Gupta, 629 and Vasudeva Varma. 2017. Deep learning for hate 630 speech detection in tweets. In Proceedings of the 631 26th international conference on World Wide Web 632 companion, pages 759-760. 633 Valerio Basile, Cristina Bosco, Elisabetta Fersini, 634 Debora Nozza, Viviana Patti, Francisco Manuel 635 Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 636 2019. Semeval-2019 task 5: Multilingual detection 637 of hate speech against immigrants and women in twit-638 ter. In Proceedings of the 13th International Work-639 shop on Semantic Evaluation, pages 54-63. ACL. 640 Tazeek Bin Abdur Rakib and Lay-Ki Soon. 2018. Using 641 the reddit corpus for cyberbully detection. pages 180-642 189. Springer. 643 Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum 644 Thain, and Lucy Vasserman. 2019. Nuanced metrics 645 for measuring unintended bias with real data for text 646 classification. In Companion proceedings of the 2019 647 world wide web conference, pages 491–500. 648 Thomas Davidson, Debasmita Bhattacharya, and Ing-649 mar Weber. 2019. Racial bias in hate speech and 650 abusive language detection datasets. arXiv preprint 651 arXiv:1905.12516. 652 Thomas Davidson, Dana Warmsley, Michael Macy, and 653 Ingmar Weber. 2017. Automated hate speech detec-654 tion and the problem of offensive language. Proceed-655 ings of the ICWSM 2017, 11(1):512–515. 656 Ona de Gibert, Naiara Perez, Aitor García-Pablos, and 657 Montse Cuadros. 2018. Hate speech dataset from a 658 white supremacy forum. In Proceedings of the ALW2 659 2018. ACL. 660 Junwei Deng and Jiaqi Ma. 2024. Computational copy-661 right: Towards a royalty model for AI music gener-662 ation platforms. In Workshop on Data-centric ma-663 chine learning. 664

773

774

- Shrey Desai and Greg Durrett. 2020. Calibration of pre-trained transformers. *arXiv preprint arXiv:2003.07892*.
- DMLR. 2024. Call for papers DMLR 2024. https: //dmlr.ai/cfp-icml24/. Accessed: 2025-02-15.
- Ann-Kathrin Dombrowski, Maximilian Alber, Christopher J Anders, Marcel Ackermann, Klaus-Robert Müller, and Pan Kessel. 2019. Explanations can be manipulated and geometry is to blame. *arXiv preprint arXiv:1906.07983*.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 214–226. ACM.

675

695

697

701

703

704

706

709

710

712

713

714

715

716

717

- Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and removing disparate impact. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 259–268. ACM.
- Antigoni-Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, others, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. *Proceedings of the ICWSM 2018*, 12(1).
- Amirata Ghorbani, Abubakar Abid, and James Zou.
  2019. Interpretation of neural networks is fragile.
  In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3681–3688.
- Lukas Grimminger and Roman Klinger. 2021. Hate towards the political opponent: A twitter corpus study of the 2020 us elections on the basis of offensive speech and stance detection. In *Proceedings of the WASSA 2021*, pages 171–180. ACL.
- Sarthak Gupta, Pranav Priyadarshi, and Manish Gupta. 2023. Hateful comment detection and hate target type prediction for video comments. In *Proceedings of the CIKM 2023*, CIKM '23, pages 3923–3927. ACM.
- Moritz Hardt, Eric Price, Nati Srebro, et al. 2016. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, pages 3315–3323.
- Peter Hase, Harry Xie, and Mohit Bansal. 2021. The out-of-distribution problem in explainability and search methods for feature importance explanations. *Advances in neural information processing systems*, 34:3650–3666.
- Marijn Janssen, Paul Brous, Elsa Estevez, Luis S Barbosa, and Tomasz Janowski. 2020. Data governance: Organizing data for trustworthy artificial intelligence. *Government information quarterly*, 37(3):101493.

- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. 2019. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*.
- Brendan Kennedy, Mohammad Atari, Aida Mostafazadeh Davani, Leigh Yeh, Ali Omrani, Yehsong Kim, others, and Morteza Dehghani. 2022. Introducing the gab hate corpus: defining and applying hate-based rhetoric to social media posts at scale. *Language Resources and Evaluation*, 56(1):79–108.
- Chris J. Kennedy, Geoff Bacon, Alexander Sahn, and Caterina von Vacano. 2020. Constructing interval variables via faceted rasch measurement and multitask deep learning: a hate speech application.
- Ravin Kohli, Matthias Feurer, Katharina Eggensperger, Bernd Bischl, and Frank Hutter. 2024. Towards quantifying the effect of datasets for benchmarking: A look at tabular machine learning. In *Workshop on Data-centric machine learning*.
- Jana Kurrek, Haji Mohammad Saleem, and Derek Ruths. 2020. Towards a comprehensive taxonomy and largescale annotated corpus for online slur usage. In *Proceedings of the WOAH 2020*, pages 138–149, Online. ACL.
- Eunjin Lee, David Braines, Mitchell Stiffler, Adam Hudler, and Daniel Harborne. 2019. Developing the sensitivity of lime for better machine learning explanation. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006, page 1100610. International Society for Optics and Photonics.
- Ilya Loshchilov, Frank Hutter, et al. 2017. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*, 5.
- Scott Lundberg. 2017. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2020. Hatexplain: A benchmark dataset for explainable hate speech detection. In *Proceedings of the AAAI 2020*.
- Ioannis Mollas, Zoe Chrysopoulou, Stamatis Karlos, and Grigorios Tsoumakas. 2022. Ethos: a multi-label hate speech detection dataset. *Complex & Intelligent Systems*, 8(6):4663–4678.
- Rafael Mosquera Gómez, Julian Eusse, Juan Ciro, Daniel Galvez, Ryan Hileman, Kurt Bollacker, and David Kanter. 2023. Speech Wikimedia: A 77 language multilingual speech dataset. In *Workshop on Data-centric machine learning*.
- Maximilian Noppel and Christian Wressnegger. 2024. Sok: Explainable machine learning in adversarial environments. In 2024 IEEE Symposium on Security and Privacy (SP), pages 2441–2459. IEEE.

- 775
- 790 791 792
- 794
- 795 796
- 797 798

- 810
- 812 813 814

811

816 817

- 818
- 819

825

829

- Nedjma Ousidhoum, Zizheng Lin, Hongming Zhang, Yangqiu Song, and Dit-Yan Yeung. 2019. Multilingual and multi-aspect hate speech analysis. In Proceedings of the EMNLP-IJCNLP 2019, pages 4675-
- Amandalynne Paullada, Inioluwa Deborah Raji, Emily M Bender, Emily Denton, and Alex Hanna. 2021. Data and its (dis) contents: A survey of dataset development and use in machine learning research. Patterns, 2(11).

4684. ACL.

- John Pavlopoulos, Jeffrey Sorensen, Lucas Dixon, Nithum Thain, and Ion Androutsopoulos. 2020. Toxicity detection: Does context really matter?
  - John Pavlopoulos, Jeffrey Sorensen, Léa Laugier, and Ion Androutsopoulos. 2021. Semeval-2021 task 5: Toxic spans detection. In Proceedings of the SemEval 2021, pages 59-69. ACL.
  - Aabha Pingle, Aditya Vyawahare, Isha Joshi, Rahul Tangsali, and Raviraj Joshi. 2023. L3CubeMahaSent-MD: A multi-domain Marathi sentiment analysis dataset and transformer models. In Workshop on Data-centric machine learning.
- Paloma Piot, Patricia Martín-Rodilla, and Javier Parapar. 2024. Metahate: A dataset for unifying efforts on hate speech detection. In Proceedings of the International AAAI Conference on Web and Social Media, volume 18, pages 2025-2039.
- 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 Proceedings of the EMNLP-IJCNLP 2019, pages 4755-4764. ACL.
- José Ribeiro, Raíssa Silva, Lucas Cardoso, and Ronnie Alves. 2021. Does dataset complexity matters for model explainers? In 2021 IEEE International Conference on Big Data (Big Data), pages 5257-5265. IEEE.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " why should i trust you?" explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining, pages 1135-1144.
- Stephen Robertson. 2004. Understanding inverse document frequency: on theoretical arguments for idf. Journal of documentation, 60(5):503-520.
- Andrea Romei and Salvatore Ruggieri. 2014. A multidisciplinary survey on discrimination analysis. The Knowledge Engineering Review, 29(05):582–638.
- Pranav Sachdeva, Ricardo Barreto, Geoff Bacon, Alexander Sahn, Caterina von Vacano, and Chris Kennedy. 2022. The measuring hate speech corpus: Leveraging rasch measurement theory for data perspectivism. In Proceedings of the LREC 2022, pages 83-94. ELRA.

Joni Salminen, Hind Almerekhi, Milos Milenkovic, Soon-gyo Jung, Jisun An, Haewoon Kwak, and Bernard Jansen. 2018. Anatomy of online hate: Developing a taxonomy and machine learning models for identifying and classifying hate in online news media. Proceedings of the ICWSM 2018, 12(1).

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

880

- Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. 2021. "everyone wants to do the model work, not the data work": Data cascades in high-stakes ai. In proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, pages 1–15.
- V Sanh. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. arXiv preprint arXiv:1910.01108.
- Rajat Shinde, Sujit Roy, Christopher E Phillips, Aman Gupta, Aditi Sheshadri, Manil Maskey, and Rahul Ramachandran. 2024. WINDSET: Weather insights and novel data for systematic evaluation and testing. In Workshop on Data-centric machine learning.
- Samuel Sithakoul, Sara Meftah, and Clément Feutry. 2024. Beexai: Benchmark to evaluate explainable ai. In World Conference on Explainable Artificial Intelligence, pages 445–468. Springer.
- Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. 2020. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In Conference on Artificial Intelligence, Ethics, and Society (AIES).
- Antonio Torralba and Alexei A Efros. 2011. Unbiased look at dataset bias. In CVPR 2011, pages 1521-1528. IEEE.
- Bertie Vidgen, Dong Nguyen, Helen Margetts, Patricia Rossini, and Rebekah Tromble. 2021. Introducing cad: the contextual abuse dataset. In Proceedings of the NAACL 2021, pages 2289–2303. ACL.
- Artem Vysogorets and Julia Kempe. 2024. Towards robust data pruning. In Workshop on Data-centric machine learning.
- T Wolf. 2019. Huggingface's transformers: State-ofthe-art natural language processing. arXiv preprint arXiv:1910.03771.
- Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2016. Ex machina: Personal attacks seen at scale.
- Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. 2024. Generalized out-of-distribution detection: A survey. International Journal of Computer Vision, pages 1-28.
- Rubing Yang, Jialin Mao, and Pratik Chaudhari. 2022. Does the data induce capacity control in deep learning? In International Conference on Machine Learning, pages 25166-25197. PMLR.

- Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, page 100211.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar.
  2019. Predicting the type and target of offensive posts in social media. In *Proceedings of the NAACL* 2019, pages 1415–1420. ACL.
- Dorothy Zhao, Alice Xiang, Jerone T A Andrews, and Orestis Papakyriakopoulos. 2024. Measuring diversity in datasets. In *Workshop on Data-centric machine learning*.
- Indre Zliobaite. 2015. A survey on measuring indirect discrimination in machine learning. *arXiv preprint arXiv:1511.00148*.

# 9 NLP Checklist

 **A1.** Did you describe the limitations of your work? Yes. Please refer to Section 8.

**A2.** Did you discuss any potential risks of your work? Yes. We have discussed some ethical considerations in Section 8.1.

**B.** Did you use or create scientific artifacts? Yes, we used existing scientific artifacts (datasets, pre-trained models, evaluation metrics).

B1. Did you cite the creators of artifacts you used?Yes. Please refer to Section 3.

B2. Did you discuss the license or terms for use
and/or distribution of any artifacts? Yes. Please
refer to Section 3.

Did you discuss if your use of existing **B3**. artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)? Our use of the MetaHate dataset follows its intended research purpose, accessed through a signed terms of use agreement. Any derivatives from our work maintain the original research-only restrictions and cannot be used outside research contexts. 

B4. Did you discuss the steps taken to check
whether the data that was collected/used contains
any information that names or uniquely identifies
individual people or offensive content, and the
steps taken to protect/anonymize it? The datasets
used contain offensive language. The sources
are publicly available, however, to avoid any
distressing feeling to our readers we avoided

presenting and cite content in the full body of the paper that can affect the readers.

**B5.** Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.? The datasets are explained in detail by the authors of the MetaHate paper. The descriptions in this paper include only what is necessary to this work.

**B6.** Did you report relevant statistics like the number of examples, details of train/test/dev splits, etc. for the data that you used/created? Yes, please refer to Sections 3 and 4.

C. Did you run computational experiments? Yes.

**C1.** Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used? Yes. Please refer to Section 3.

**C2.** Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? Yes. We did not perform hyperparameter tuning.

**C3.** Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run? Yes. We report performance results in Section 4 and Appendix A. **C4.** If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation, such as NLTK, Spacy, ROUGE, etc.), did you report the implementation, model, and parameter settings used? Yes. Please refer to Section 3.

**D.** Did you use human annotators (e.g., crowd-workers) or research with human participants? No. **D1.** Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.? N/A.

**D2.** Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)? N/A.

**D3.** Did you discuss whether and how consent was obtained from people whose data you're using/curating? N/A.

**D4.** Was the data collection protocol approved (or determined exempt) by an ethics review board? N/A.

**D5.** Did you report the basic demographic and geographic characteristics of the annotator

- population that is the source of the data? N/A. 986 E. Did you use AI assistants (e.g., ChatGPT, Copi-987 lot) in your research, coding, or writing? Yes. We 988 used AI language models for proofreading portions 989 of the manuscript to check for grammatical errors 990 991 and clarity. E1. Did you include information about your use of 992 AI assistants? Yes. Please refer to Section 7. 993
- 994

# A Performance comparison of LR and DB across different hate speech datasets

Dataset	<b>F1</b>		AUROC		Precision		Recall		Accuracy		Bal. Acc	
	LR	DB	LR	DB	LR	DB	LR	DB	LR	DB	LR	DB
MLMA 2019	0.433	0.433	0.534	0.512	0.000	0.000	0.000	0.000	0.763	0.765	0.499	0.500
HatEval 2019	0.724	0.760	0.819	0.854	0.719	0.734	0.610	0.696	0.739	0.769	0.720	0.759
News Media 2018	0.848	0.886	0.938	0.960	0.937	0.941	0.878	0.928	0.871	0.907	0.865	0.890
Hate Speech 2020-22	0.698	0.744	0.844	0.870	0.718	0.716	0.407	0.523	0.802	0.820	0.675	0.725
Offensive 2017	0.566	0.537	0.881	0.875	0.684	0.640	0.091	0.056	0.945	0.944	0.544	0.527
Toxic Spans 2021	0.479	0.479	0.596	0.623	0.918	0.918	1.000	1.000	0.918	0.918	0.500	0.500
CAD 2021	0.574	0.691	0.769	0.830	0.791	0.758	0.144	0.336	0.831	0.854	0.568	0.656
HateComments 2023	0.725	0.782	0.830	0.856	0.688	0.805	0.759	0.718	0.725	0.785	0.727	0.781
Supremacist 2018	0.535	0.637	0.842	0.897	0.727	0.762	0.069	0.207	0.895	0.906	0.533	0.599
Slur Corpus 2020	0.805	0.882	0.880	0.946	0.806	0.894	0.816	0.874	0.805	0.882	0.805	0.882
HateXplain 2020	0.792	0.795	0.860	0.881	0.772	0.859	0.739	0.653	0.798	0.809	0.790	0.787
ExMachina 2016	0.826	0.868	0.950	0.973	0.880	0.920	0.568	0.657	0.932	0.947	0.778	0.824
Context Toxic 2020	0.497	0.497	0.678	0.711	0.000	0.000	0.000	0.000	0.988	0.988	0.500	0.500
ENCASE 2018	0.920	0.927	0.964	0.975	0.891	0.886	0.875	0.905	0.937	0.942	0.918	0.930
ETHOS 2022	0.624	0.736	0.693	0.837	0.500	0.620	0.515	0.721	0.660	0.755	0.625	0.747
US Elections 2020	0.469	0.762	0.667	0.922	0.000	0.811	0.000	0.435	0.885	0.923	0.500	0.711
Jigsaw Toxic	0.767	0.835	0.959	0.975	0.857	0.743	0.408	0.641	0.962	0.967	0.702	0.814
OLID 2019	0.672	0.764	0.774	0.860	0.775	0.756	0.378	0.604	0.753	0.800	0.661	0.752
BullyDetect 2018	0.758	0.839	0.893	0.940	0.855	0.736	0.489	0.815	0.834	0.867	0.728	0.851
Gab Hate 2022	0.559	0.666	0.847	0.885	0.737	0.688	0.090	0.255	0.920	0.927	0.543	0.622
Intervene 2019	0.888	0.904	0.930	0.943	0.909	0.899	0.825	0.874	0.893	0.907	0.883	0.902

995

# **B** Impact of source training data on features attribution.

Figure 8: SHAP explanation of the LR model trained on bullydetect\_2018 (above) and ENCASE\_2018 (below) and tested on the same out of distribution sentence. The scores relate to the predicted probability of the positive class (hate), namely, PP = 1 (above) and PP = 0.99 (below)



makes you cry from the pain every time you shift you fucking subhuman scum.....Fuck you. I would gladly cut your fucking balls off in a minute....go suck cock you piece of shift Fuck you and your mother you sorry fuckin swill

Figure 9: SHAP explanation of the LR model trained on OLID\_2019 (above) and ethos\_2022 (below) and tested on the same out of distribution sentence. The scores relate to the predicted probability of the positive class (hate), namely, PP = 0.96 (above) and PP = 0.94 (below)

inputs
The righteous ALT LEFT will eradicate fascist inbred redneck white supremacisk like cockroaches. We will hunt them down like pigs and lynch them in our town squares
and cities. We will gut them as they emerge from their backwoods swamps. Be warned white supremacist/nationalist scum. We will fullyou and burn down your homes and
destroy you and your families. The righteous ALT LEFT will cut your filthy throats in your homes in fort of your children

inputs
The righteous ALT LEFT will eradicate **fascisti** inbred redneck white supremacists like cockroaches. We will hunt them down like pigs and lynch them in our town squares
and cities. We will gut them as they emerge from their backwoods swamps. Be warned white supremacist/nationalist scum. We will kill you and burn down your homes and
destroy you and your families. The righteous ALT LEFT will cut your filthy throats in your homes in fort of your children