

Robustness Evaluation of Hate Speech Detection Models Under Structured Adversarial Obfuscation

Anonymous ACL submission

Abstract

Hate speech classifiers are vulnerable to adversarial text obfuscation, yet existing robustness benchmarks typically evaluate models using either random noise or isolated perturbation techniques, failing to capture the structured multi-strategy evasion behaviour observed in practice. We present a systematic benchmark evaluating four transformer-based detectors against six obfuscation techniques under three attack regimes spanning *realistic*, *uniform*, and *adversarial* distributions. Beyond prior work, we additionally model realistic social media language by obfuscating neutral (non-toxic) tokens alongside toxic ones, and introduce *obfuscation intensity* as a dedicated evaluation axis. Our results show substantial robustness degradation across all models, with F1 drops of up to -0.356 and performance deteriorating non-linearly as perturbation density increases. Obfuscation-aware fine-tuning recovers up to $+0.326$ F1, demonstrating that robustness can be substantially improved without architectural modification.

1 Introduction

The proliferation of hate speech on social media platforms represents one of the most pressing challenges in online safety. A recent cross-national survey found that two in three people frequently encounter hate speech online (Ipsos, 2023), while platforms such as Facebook and Instagram removed over 14 million pieces of hateful content in a single quarter of 2024 alone (Meta Platforms, 2024). Beyond sheer volume, hate speech causes measurable psychological harm to targeted individuals and has been linked to real-world violence and discrimination (Vidgen and Derczynski, 2020). Automated detection systems have thus become an indispensable tool for large-scale content moderation.

The advent of transformer-based models has substantially advanced the state of the art in hate

speech detection (Waseem and Hovy, 2016; Davidson et al., 2017; Caselli et al., 2021). Models such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and domain-adapted variants like HateBERT (Caselli et al., 2021) achieve strong performance on standard benchmarks. However, these systems are evaluated almost exclusively on clean, unperturbed corpora — an assumption that rarely holds in adversarial real-world conditions.

In practice, users attempting to evade automated moderation employ *textual obfuscation*: deliberate manipulation of word surface forms to preserve semantic intent while bypassing classifiers. Strategies range from character substitution (*bitch* \rightarrow *b!tch*) and phonetic replacement (*fool* \rightarrow *phool*) to intra-word spacing (*moron* \rightarrow *m o r o n*) and composite perturbations. Such evasion is well documented in the content moderation literature (Gröndahl et al., 2018), yet its systematic impact on transformer-based detectors remains critically understudied.

Existing work on obfuscation robustness suffers from two key limitations. Prior benchmarks either apply uniform random noise to all tokens — overestimating evasion impact — or target single isolated words, underestimating the structured multi-strategy behaviour of real adversaries (Aggarwal and Zesch, 2022). Moreover, evaluations treat obfuscation as binary, ignoring how perturbation *density* affects detection difficulty. We address both gaps with a comprehensive benchmark that models obfuscation as a probabilistic, policy-driven process spanning real-world, uniform, and adversarial attack distributions. We further introduce obfuscation intensity as a novel evaluation axis and construct a realistic dataset variant that obfuscates neutral tokens alongside toxic ones, faithfully simulating the informal register of social media users who casually abbreviate benign language while concealing harmful content.

Specifically, we examine: (1) the extent to which obfuscation degrades detection performance across

084 model architectures, (2) whether certain transform- 131
085 ers exhibit greater robustness to adversarial per- 132
086 turbations, (3) how obfuscation intensity affects 133
087 model performance, and (4) whether obfuscation- 134
088 aware training can recover lost robustness. 135

089 Contributions. 136

- 090 • We present the **first systematic benchmark** 137
091 evaluating four transformer models against 138
092 six obfuscation techniques under three proba- 139
093 bilistic policy regimes, with F1 drops of up to 140
094 -0.356 under realistic evasion conditions. 141
- 095 • We introduce **obfuscation intensity** as a novel 142
096 evaluation axis and a **realistic dataset vari-** 143
097 **ant** extending toxic-word perturbation to neu- 144
098 tral tokens, exposing a non-linear degradation 145
099 curve absent from prior work. 146
- 100 • We show that **obfuscation-aware fine-tuning** 147
101 recovers up to $+0.326$ F1 with no architectural 148
102 changes, providing a practical deployable de- 149
103 fence for content moderation systems. 150

104 2 Related Work 151

105 2.1 Hate Speech Detection 152

106 Early approaches to hate speech detection relied 153
107 on lexicon-based methods and traditional machine 154
108 learning classifiers such as SVMs and logistic re- 155
109 gression (Waseem and Hovy, 2016; Davidson et al., 156
110 2017). Davidson et al. (2017) introduced the widely 157
111 adopted three-class Twitter corpus distinguishing 158
112 hate speech, offensive language, and neither, which 159
113 remains a standard benchmark. The introduction 160
114 of transformer architectures (Devlin et al., 2019) 161
115 marked a step change in detection performance. 162
116 Caselli et al. (2021) demonstrated that domain- 163
117 adaptive pre-training on abusive online content 164
118 yields substantial gains over general-purpose en- 165
119 coders, motivating the inclusion of HateBERT in 166
120 our benchmark. Fortuna and Nunes (2018) pro- 167
121 vided a comprehensive survey of hate speech detec- 168
122 tion methods, highlighting the persistent challenge 169
123 of generalisation across datasets and linguistic re- 170
124 gisters — a challenge that obfuscation directly ex- 171
125 acerbates. 172

126 2.2 Adversarial Robustness in NLP 173

127 The vulnerability of NLP models to adversarial per- 174
128 turbations has been extensively studied. Ebrahimi 175
129 et al. (2018) introduced HotFlip, a white-box at- 176
130 tack that exploits gradients to flip characters or 177

tokens to change classification outcomes. Wallace 131
et al. (2019) demonstrated that a fixed sequence of 132
tokens, optimised by gradient search, could consis- 133
tently reduce accuracy across many inputs, estab- 134
lishing the concept of universal adversarial triggers. 135
Morris et al. (2020) unified this literature through 136
TextAttack, a modular framework for adversarial 137
attacks and data augmentation in NLP, providing 138
standardised benchmarking infrastructure. Collec- 139
tively, these works establish that surface-level per- 140
turbations — even character-level ones — can fun- 141
damentally undermine model robustness, a finding 142
we extend specifically to the hate speech detection 143
domain. 144

145 2.3 Obfuscation as Evasion of Hate Speech 146 147 Detection 148

149 The intersection of adversarial NLP and hate 150
151 speech moderation has received growing attention. 152
Gröndahl et al. (2018) conducted a seminal study 153
demonstrating that all proposed hate speech de- 154
tection techniques are brittle against adversaries 155
who insert typos, change word boundaries, or add 156
innocuous words. Crucially, they showed that 157
character-level features make models more attack- 158
resistant than word-level features — an insight that 159
informs our choice of subword tokeniser-based 160
models. Röttger et al. (2021) introduced Hate- 161
Check, a suite of functional tests for hate speech 162
detection models, revealing critical weaknesses in 163
both academic and commercial systems. While 164
HateCheck evaluates models on handcrafted test 165
cases, it does not systematically vary obfuscation 166
type or intensity across a realistic corpus, a gap our 167
benchmark addresses directly. 168

169 Most closely related to our work, Aggarwal and 170
171 Zesch (2022) analysed the real vulnerability of hate 172
173 speech detection systems against targeted inten- 174
175 tional noise, conducting a user study to identify 176
177 which words users would actually obfuscate in a 178
179 post. They found that real-world vulnerability is al- 179
180 most as high as white-box attacks, and much more 180
181 severe than non-targeted dictionary methods. Our 181
work differs in three important ways. First, rather 182
than constraining obfuscation to a single token per 183
post, we apply structured multi-technique perturba- 184
tion under three probabilistic policy regimes, mod- 185
elling a wider spectrum of adversarial behaviour. 186
Second, we introduce obfuscation intensity as an 187
explicit evaluation axis, quantifying how detection 188
difficulty scales with perturbation density. Third, 189
we extend evaluation beyond single-model settings 190

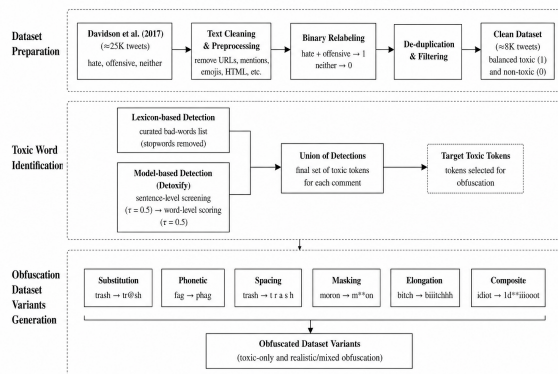


Figure 1: Overview of the obfuscation pipeline: dataset preparation (top), hybrid toxic word identification (middle), and obfuscation dataset variant generation (bottom).

to a controlled four-model benchmark, enabling cross-architecture robustness comparison that is absent from prior work. Existing work largely evaluates isolated perturbations or handcrafted attacks, but lacks a unified benchmark framework that jointly studies obfuscation type, perturbation density, and robustness recovery under controlled conditions.

3 Dataset and Obfuscation Framework

3.1 Source Dataset and Label Schema

We build upon the benchmark dataset introduced by Davidson et al. (2017), which comprises approximately 25,000 tweets annotated into three categories: *hate speech*, *offensive language*, and *neither*. For our binary classification objective, we consolidate the first two categories into a single positive class (label 1) representing toxic or hate speech, while retaining the third as the negative class (label 0). This relabelling is consistent with prior work that treats offensive and hate speech as a unified threat class when evaluating content moderation systems (Fortuna and Nunes, 2018; Founta et al., 2018). The overall dataset preparation and obfuscation generation pipeline is illustrated in Figure 1.

The final working corpus comprises 8,326 samples with a balanced class distribution of 4,163 toxic and 4,163 non-toxic instances, constructed through random sampling from the merged toxic class. This balanced setup ensures that accuracy and F1 are directly comparable across all experimental conditions without class-imbalance confounds.

3.2 Toxic Word Identification

A prerequisite for targeted obfuscation is accurate identification of toxic tokens within each utterance. Naively obfuscating all tokens would overestimate real-world evasion impact (Aggarwal and Zesch, 2022), while random single-token selection would underestimate it. We therefore employ a *hybrid detection strategy* that combines two complementary signals, as illustrated in the middle panel of Figure 1.

A lexicon-based filter matches surface forms against a curated bad-words list,¹ excluding common stopwords to suppress false positives; while effective for known slurs, it fails to capture context-dependent toxic language absent from the vocabulary (Fortuna and Nunes, 2018). To address this, we deploy *Detoxify* (Han and Unitary Team, 2020), which first screens the full utterance at sentence level ($\tau = 0.5$), then scores each candidate token individually (length > 3 , non-stopword), retaining those exceeding τ — a threshold widely adopted in toxicity pipelines to balance precision and recall (Vidgen and Derczynski, 2020). The final toxic token set is the *union* of both detectors, ensuring high recall of context-dependent slurs that evade purely lexical methods while preserving coverage of well-known offensive terms that may fall below the model threshold in isolation.

3.3 Obfuscation Techniques

We implement six character- and token-level obfuscation operators, each simulating a distinct evasion strategy documented in real-world hate speech evasion (Gröndahl et al., 2018; Aggarwal and Zesch, 2022). The full pipeline — from dataset preparation through toxic word identification to obfuscation variant generation — is illustrated in Figure 1. Table 1 summarises each technique with a representative example.

Substitution applies a probabilistic character mapping ($p = 0.6$ per eligible character) using a leet-speak substitution table (e.g., $a \leftrightarrow @$, $e \leftrightarrow 3$, $i \leftrightarrow !$), guaranteeing at least one substitution per toxic token. **Phonetic** obfuscation replaces substrings with phonetically equivalent sequences (e.g., $ph \leftrightarrow f$, $oo \leftrightarrow u$) while protecting common digraphs (ch , sh , th) to preserve readability. **Spacing** inserts intra-word whitespace via single-split, double-split,

¹Bad-words list sourced from Carnegie Mellon University: <https://www.cs.cmu.edu/~biglou/resources/bad-words.txt>

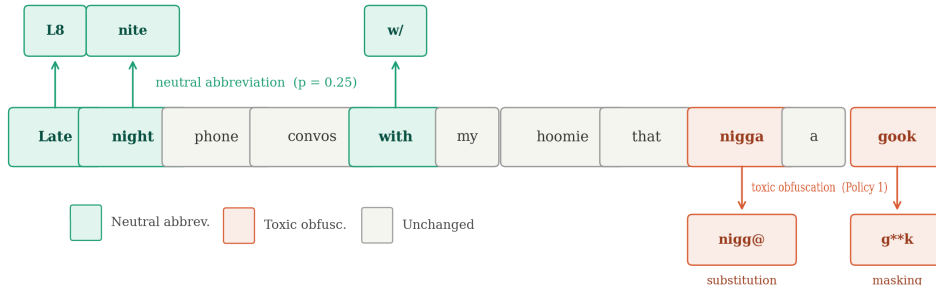


Figure 2: Illustration of the two-stage obfuscation pipeline on a representative comment from D_c . Neutral tokens are stochastically abbreviated upward ($p = 0.25$); toxic tokens are obfuscated downward via the Policy 1 selected technique (substitution: *nigga* \rightarrow *nigg@*; masking: *gook* \rightarrow *g**k*). Unchanged tokens are shown in grey.

Technique	Example
Substitution	trash \rightarrow tr@sh
Phonetic	fool \rightarrow phool
Spacing	hate \rightarrow h a t e
Masking	moron \rightarrow m**on
Elongation	bitch \rightarrow biiitchhh
Composite	idiot \rightarrow 1d**iiiioot

Table 1: Obfuscation techniques used for generating adversarial toxic text variants.

Dataset	Size	Description
D_a (Clean)	8,326	Original corpus; no perturbation
D_b (Mixed)	8,326	Toxic tokens obfuscated via 3-policy blend
D_c (Realistic)	8,326	D_b + neutral word abbreviation ($p=0.25$)
D_d (Individual)	8,326 each	One technique per variant (D_1 – D_6)

Table 2: Dataset variants used for benchmarking.

or character-level decomposition, directly exploiting the token boundary assumptions of subword tokenisers such as BPE (Sennrich et al., 2016). **Masking** replaces a contiguous interior substring with special symbols ($*$, $\#$), obscuring the word’s core characters while preserving its overall shape. **Elongation** stochastically repeats characters 2–4 times at randomly selected positions, mimicking the emphatic spelling common in informal social media text (Eisenstein, 2013). Finally, **Composite** sequentially chains 2–3 randomly sampled techniques from the above set, generating the most adversarially challenging surface forms and most closely simulating sophisticated evasion behaviour observed in practice (Aggarwal and Zesch, 2022).

3.4 Dataset Variants

We construct four dataset variants (Table 2) to enable controlled benchmarking across a spectrum of obfuscation complexity. Figure 2 illustrates the two-stage obfuscation process on a representative

comment.

D_a (**Baseline**) serves as the unperturbed reference, retaining the original cleaned text with no artificial perturbation. D_b (**Mixed Obfuscation**) introduces policy-level stochasticity: for each comment, one of three calibrated probability distributions (Table 3) is selected with equal probability ($P = 0.33$), and for each toxic token the technique is sampled from the chosen policy. D_c (**Realistic**) extends D_b by additionally applying an internet-abbreviation lexicon to non-toxic tokens ($P = 0.25$ per entry), simulating the informal register prevalent in social media and mitigating the risk that models learn to associate obfuscation patterns with toxicity rather than semantic content (Eisenstein, 2013; Tagliamonte and Denis, 2008). D_d (**Individual**) consists of six independent variants (D_1 – D_6), each applying a single obfuscation technique exclusively to identified toxic tokens, enabling fine-grained attribution of model degradation to individual perturbation types.

3.5 Policy Distributions for Mixed Obfuscation

Three policies govern technique selection within D_b and D_c (Table 3). Unlike prior synthetic perturbation frameworks that apply uniformly random transformations (Gröndahl et al., 2018; Jones et al., 2020), our policies collectively span the space of plausible real-world evasion behaviours through non-uniform obfuscation priors.

P1 (Real-world) assigns higher probability to substitution (0.30) and phonetic alteration (0.20), consistent with prior observations that character-level manipulation is the most prevalent evasion strategy in online hate speech corpora (Gröndahl et al., 2018). **P2 (Uniform)** assigns equal probability (≈ 0.17) to all techniques, eliminating distributional bias and serving as a bias-free experimental

Policy	Subst.	Phonetic	Spacing	Masking	Elong.	Composite
P1 (Real-world)	0.30	0.20	0.10	0.05	0.15	0.20
P2 (Uniform)	0.17	0.17	0.17	0.17	0.17	0.17
P3 (Adversarial)	0.10	0.10	0.25	0.15	0.10	0.30

Table 3: Sampling distributions over obfuscation techniques for the three policy regimes used in mixed-obfuscation generation, modelling realistic, uniform, and adversarial evasion behaviour.

control for technique-level sensitivity analysis. **P3 (Adversarial)** up-weights composite (0.30) and spacing (0.25), modelling a worst-case attacker who preferentially selects transformations that fragment subword tokenisation boundaries — the primary structural assumption of transformer-based encoders (Ebrahimi et al., 2018; Wallace et al., 2019). For each comment, a single policy is drawn uniformly ($P = 0.33$) and applied consistently to all toxic tokens, preserving intra-comment stylistic coherence.

3.6 Neutral Word Obfuscation in D_c

To construct a maximally realistic evaluation set, D_c augments the mixed-obfuscation layer with stochastic substitution of neutral tokens using a curated lexicon of 80+ internet abbreviations, compiled from established internet language resources (Han and Baldwin, 2011; Eisenstein, 2013) and spanning four categories reflective of authentic social media register: chat shortcuts (*great*→*gr8*), internet slang (*oh my god*→*omg*), vowel dropping (*people*→*ppl*), and casual reductions (*kind of*→*kinda*). This categorisation is grounded in prior sociolinguistic work documenting the systematic nature of informal language on social media platforms (Tagliamonte and Denis, 2008).

Multi-word entries are matched with longest-first priority to prevent partial overlaps, and original capitalisation is preserved via case-aware replacement. Crucially, toxic tokens identified in §3.2 are explicitly protected from neutral obfuscation, preventing conflation of the two transformation layers. Each neutral token is substituted independently with probability $p = 0.25$, ensuring naturalness without over-abbreviation. As illustrated in Figure 2, the result is a comment where both harmful and benign tokens reflect realistic social media writing patterns simultaneously.

4 Experimental Setup

4.1 Models

We evaluate four transformer-based encoder models that collectively span the space of lightweight,

Model	HuggingFace ID	Params	Pre-training
DistilBERT	distilbert-base-uncased	66M	General (EN)
RoBERTa	roberta-base	125M	General (EN)
HateBERT	GroNLP/hateBERT	110M	Hate speech
XLNet	xlnet-base	270M	Multilingual

Table 4: Pre-trained encoder models with HuggingFace identifiers, parameter counts, and pre-training corpora.

high-capacity, domain-specialised, and multilingual architectures (Table 4). This selection is motivated by the goal of isolating whether obfuscation robustness is a function of model scale, pre-training data breadth, or domain specialisation — three axes identified as critical determinants of robustness in adversarial NLP (Morris et al., 2020; Wang et al., 2021).

Table 4 summarises the four models. DistilBERT (Sanh et al., 2019) represents efficiency-constrained deployments; RoBERTa (Liu et al., 2019) establishes the general-purpose ceiling; HateBERT (?) tests whether domain specialisation confers robustness; and XLNet (Conneau et al., 2020) evaluates cross-lingual pre-training breadth. All models are sourced from the HuggingFace Model Hub (Wolf et al., 2020) and fine-tuned under identical conditions.

We deliberately exclude autoregressive large language models such as GPT-4 (OpenAI, 2023) and LLaMA (Touvron et al., 2023), as encoder-only architectures remain the dominant paradigm in deployed content moderation systems owing to their inference efficiency and direct suitability for sequence classification (Vidgen and Derczynski, 2020; Modha et al., 2020).

4.2 Experimental Conditions

We define six experimental conditions that progress from clean baseline evaluation through increasingly adversarial settings to a recovery analysis, summarised in Table 5.

CLEAN evaluates all models on unperturbed data, establishing the upper-bound performance ceiling under standard conditions. MIXED and REALISTIC transfer models trained on clean data to obfuscated test sets D_b and D_c respectively, directly

Condition	Train	Test	Objective
CLEAN	D_a	D_a	Baseline performance ceiling
MIXED	D_a	D_b	Robustness under mixed obfuscation
REALISTIC	D_a	D_c	Robustness under realistic obfuscation
PER-TECHNIQUE	D_a	$D_d (\times 6)$	Per-technique vulnerability analysis
INTENSITY	D_a	D_b (binned)	Effect of obfuscation density
RECOVERY	D_b	D_b	Obfuscation-aware training recovery

Table 5: Experimental conditions with train/test configurations and objectives.

measuring robustness degradation under controlled evasion. PER-TECHNIQUE disaggregates performance across each of the six individual obfuscation variants in D_d , enabling fine-grained attribution of degradation to specific perturbation types. INTENSITY partitions D_b by the number of obfuscated tokens per comment (0, 1, 2–3, ≥ 4), quantifying how detection difficulty scales with perturbation density — an evaluation axis absent from prior work. Finally, RECOVERY trains models directly on D_b and evaluates on D_b , testing whether obfuscation-aware fine-tuning can recover performance lost under adversarial conditions.

4.3 Training and Evaluation

All datasets are partitioned using a stratified 80/20 train–test split to preserve class balance across evaluation conditions. Full hyperparameter and implementation details are provided in Appendix A.

We adopt **F1** as the primary metric (Fortuna and Nunes, 2018; Vidgen and Derczynski, 2020), with full metrics (Accuracy, Precision, Recall, F1) reported for main conditions (Table 6) and F1 only for per-technique results (Table 7). Recall is additionally reported for intensity analysis to directly capture missed hate speech as perturbation density increases.

5 Results and Discussion

5.1 Baseline Performance and Obfuscation Impact

All models achieve strong F1 on clean data (D_a : 0.936–0.945), confirming that standard fine-tuning is sufficient for unperturbed hate speech detection. Under obfuscation (D_b), performance degrades substantially across all models (between -0.153 and -0.356 F1). **DistilBERT** suffers the largest drop (F1 = 0.587, $\Delta = -0.356$), while **RoBERTa** is most resilient (F1 = 0.792, $\Delta = -0.153$). The dominant failure mode is a **recall collapse**: Precision remains stable (>0.90) while Recall plummets (DistilBERT: 0.938 \rightarrow 0.428), indicating models miss obfuscated hate speech rather than mislabel

neutral content — a particularly dangerous failure for deployed moderation systems. This is primarily a *tokenisation problem*: substitution and spacing decompose toxic tokens into rare BPE subword fragments (Sennrich et al., 2016) that carry no toxicity signal, rendering harmful content invisible to the classifier. Elongation — the only technique preserving subword boundaries — is consistently easiest to detect (F1: 0.808–0.878), directly confirming this hypothesis. Comparing D_b and D_c ($\Delta F1 \leq 0.003$), neutral-token abbreviation contributes negligibly to degradation, confirming models respond to *semantic toxicity signals* rather than surface obfuscation patterns — a confound that prior benchmarks (Gröndahl et al., 2018) could not rule out.

5.2 Per-Technique Vulnerability

Table 7 reveals pronounced variation across techniques. **Substitution** is most effective (DistilBERT: F1 = 0.612; HateBERT: 0.657), as symbol replacement directly disrupts subword segmentation. **Masking** causes similar degradation for weaker models yet RoBERTa handles it well (F1 = 0.838), suggesting its aggressive pre-training (Liu et al., 2019) yields more robust subword representations. **Elongation** is the weakest attack (F1: 0.808–0.878) as character repetition preserves recoverable subword structure. **Composite** chaining produces intermediate degradation, demonstrating that technique combination does not simply accumulate individual effects. Notably, **HateBERT** does not outperform general-purpose encoders despite domain specialisation ($\Delta F1 = -0.292$ vs RoBERTa’s -0.153), extending Röttger et al. (2021) and Aggarwal and Zesch (2022): domain adaptation improves recognition of *known* toxic surface forms but confers no advantage against deliberate surface disruption. **Robustness requires surface invariance, not domain familiarity.**

5.3 Intensity Analysis

Figure 3 reveals a **non-monotonic** degradation pattern. The sharpest F1 drop occurs at low intensity (1 token), reflecting the *toxic anchor effect*: classifiers rely disproportionately on one or two salient toxic tokens, so corrupting even one collapses detection confidence (DistilBERT: -0.240). At medium intensity (2–3 tokens), stronger models partially recover as structural anomalies — unusual punctuation density, alternating symbol-letter sequences — provide indirect toxicity signals. At high inten-

Dataset	Metric	DistilBERT	RoBERTa	HateBERT	XML-R
D_a (Clean)	Accuracy	0.994 [†]	0.944	0.945	0.935
	Precision	0.949	0.933	0.943	0.938
	Recall	0.938	0.958	0.947	0.934
	F1	0.943	0.945	0.945	0.936
D_b (Mixed Obf.)	Accuracy	0.698	0.819	0.731	0.778
	Precision	0.932	0.933	0.921	0.914
	Recall	0.428	0.687	0.506	0.614
	F1	0.587	0.792	0.653	0.735
	$\Delta F1$ vs D_a	-0.356	-0.153	-0.292	-0.201
D_c (Realistic Obf.)	Accuracy	0.695	0.818	0.732	0.778
	Precision	0.928	0.934	0.923	0.915
	Recall	0.423	0.684	0.507	0.613
	F1	0.582	0.790	0.655	0.734
	$\Delta F1$ vs D_a	-0.361	-0.155	-0.290	-0.202

Table 6: Main results across clean and obfuscated conditions. $\Delta F1$ vs D_a shows drop from clean baseline, shaded light-to-deep red proportional to magnitude. Best F1 per dataset **bolded**. [†]DistilBERT accuracy reflects a slight positive-class prediction bias; F1 remains the primary metric.

Technique	DistilBERT	RoBERTa	HateBERT	XML-R
Substitution	0.612	0.770	0.657	0.713
Phonetic	0.766	0.765	0.774	0.779
Spacing	0.686	0.761	0.729	0.757
Masking	0.640	0.838	0.683	0.789
Elongation	0.808	0.878	0.825	0.870
Composite	0.723	0.837	0.756	0.816

Table 7: F1 scores per obfuscation technique (D_a , D1–D6). Cells shaded red-to-teal: **red** = lower F1, **teal** = higher F1.

sity (≥ 4 tokens), *signal saturation* drives further recovery (RoBERTa: F1 = 0.905), as pervasive surface corruption renders comments structurally anomalous in ways detectable independently of toxic token recognition. Critically, **moderate obfuscation (1–3 tokens) is the most dangerous evasion regime**: a sophisticated adversary maximises evasion by surgically targeting only the most salient toxic tokens, staying below the threshold at which structural anomaly detection activates. This intensity axis was entirely absent from prior benchmarks (Gröndahl et al., 2018; Aggarwal and Zesch, 2022).

5.4 Obfuscation-Aware Training

Table 8 shows that obfuscation-aware fine-tuning ($D_b \rightarrow D_b$) recovers substantial performance across all models, with all recovering to within 0.03 F1 of their clean baseline. DistilBERT achieves the

largest gain (+0.326), confirming its vulnerability stems from distributional mismatch rather than architectural limitations. Unlike gradient-based adversarial training (Ebrahimi et al., 2018; Wallace et al., 2019), our approach requires only synthetically generated obfuscated data — producible from any existing hate speech corpus without human annotation. We recommend obfuscation-aware fine-tuning as a **first-line, deployable defence** for any content moderation pipeline, particularly efficiency-constrained systems where DistilBERT’s recovery gain is largest.

6 Conclusion

This work presents a structured robustness benchmark for evaluating hate speech detection systems under adversarial text obfuscation. Across six obfuscation techniques, multiple policy regimes,

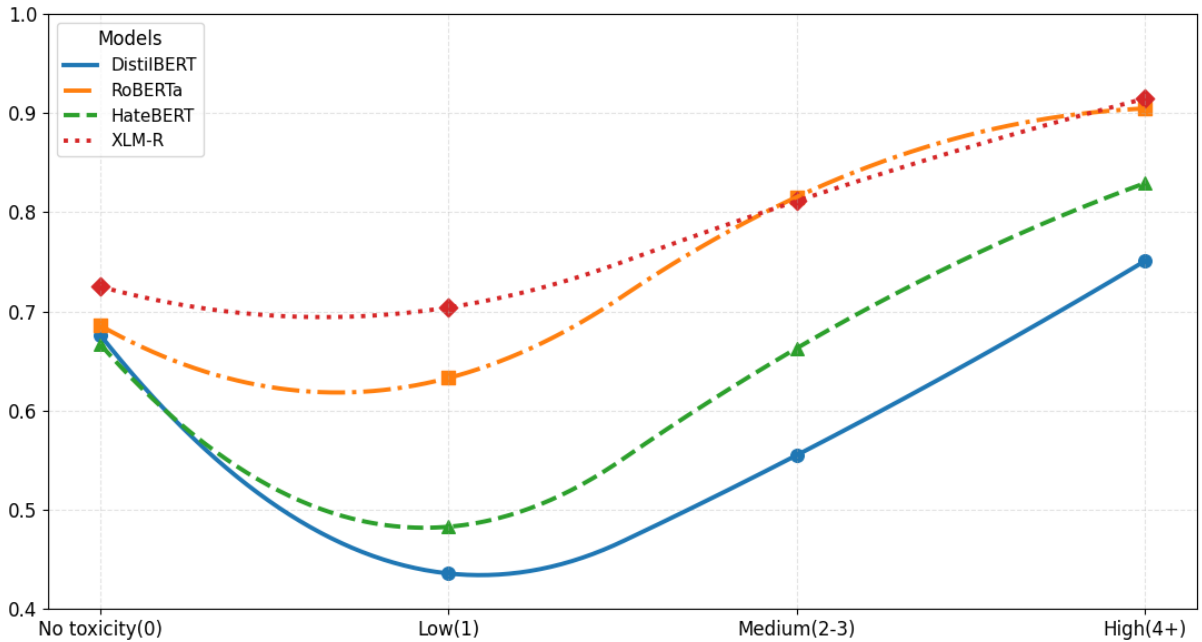


Figure 3: F1 scores across obfuscation intensity bins. Performance is non-monotonic: low intensity (1 token) causes the sharpest drop via the toxic anchor effect, with partial recovery at medium (2–3 tokens) and substantial recovery at high intensity (≥ 4 tokens) via signal saturation.

Model	$D_a \rightarrow D_a$	$D_a \rightarrow D_b$	Δ Degradation	$D_b \rightarrow D_b$	Δ Recovery
DistilBERT	0.943	0.587	-0.356	0.913	+0.326
RoBERTa	0.945	0.792	-0.153	0.928	+0.136
HateBERT	0.945	0.653	-0.292	0.913	+0.260
XLM-R	0.936	0.735	-0.201	0.905	+0.170

Table 8: F1 performance under clean, obfuscated, and obfuscation-aware training conditions. Δ Degradation and Δ Recovery indicate performance change relative to the clean baseline, with shading proportional to magnitude.

and four transformer architectures, our experiments show that modern hate speech detectors remain highly sensitive to surface-level perturbations despite strong performance on clean benchmarks. Our results demonstrate that even state-of-the-art models are highly vulnerable to deliberate surface-level manipulation, with F1 drops of up to -0.356 driven primarily by tokenisation disruption rather than semantic confusion. Critically, domain specialisation alone does not confer robustness — pre-training breadth and vocabulary coverage are stronger determinants, with RoBERTa consistently outperforming the domain-adapted HateBERT across all conditions. Our intensity analysis identifies moderate obfuscation (1–3 tokens) as the most dangerous evasion regime, a finding absent from prior work. Finally, obfuscation-aware fine-tuning recovers up to $+0.326$ F1 with no architectural changes, offering a practical, scalable, and immediately deployable defence for content mod-

eration systems. We release our dataset variants and obfuscation framework to support reproducible evaluation of future hate speech detection systems under adversarial conditions.

7 Limitations

Synthetic obfuscation. Our obfuscation datasets are algorithmically generated using probabilistic mechanisms rather than collected from real-world adversarial posts. While our policy framework is grounded in documented evasion strategies (Gröndahl et al., 2018), synthetically obfuscated text may not fully capture the creativity, context-sensitivity, and linguistic diversity of human-authored evasion.

Obfuscation coverage. Our six techniques, while well-documented, represent a subset of real-world evasion strategies. Emoji-based obfuscation (Kirk et al., 2022), code-switching, homoglyph attacks, zero-width character insertion, and LLM-

561	generated semantic paraphrases that preserve hate-	Harmful content. Our benchmark is derived	610
562	ful intent while evading surface-level detection are	from the dataset of Davidson et al. (2017) , which	611
563	not covered, and real adversaries may combine	contains real instances of hate speech and offensive	612
564	these with our studied techniques in unpredictable	language. We do not create new harmful content;	613
565	ways.	all toxic examples originate from the source corpus	614
566	Dataset and language scope. Our benchmark de-	and are used solely for research and evaluation pur-	615
567	rides from a single English Twitter corpus (David-	poses. Researchers using the dataset should handle	616
568	son et al., 2017) collected in 2017. Hate speech	it responsibly and follow institutional guidelines	617
569	evolves rapidly across platforms and languages,	for working with sensitive material.	618
570	and findings may not generalise to more recent,	Societal impact. Improving robustness in hate	619
571	multilingual, or platform-diverse settings. Obfus-	speech detection can help reduce online abuse and	620
572	cation impact on subword tokenisation may also	better protect marginalised communities. How-	621
573	differ substantially across morphologically rich or	ever, automated moderation systems also risk over-	622
574	non-Latin-script languages.	ensorship, demographic bias, or misuse in restric-	623
575	Model and defence scope. We focus on encoder-	tive surveillance settings (Vidgen and Derczynski,	624
576	only architectures; the robustness of generative	2020). We therefore encourage deployment prac-	625
577	LLMs in moderation contexts remains unstudied.	tices that include transparency, human oversight,	626
578	Our recovery evaluation also assumes a fixed ob-	and regular auditing.	627
579	fuscation distribution at test time — adaptive adver-	9 Acknowledgements	628
580	saries who observe and circumvent the defended	The authors used AI-assisted tools, including	629
581	model may require more sophisticated defence	Claude (Anthropic), ChatGPT, and Grammarly,	630
582	strategies (Wallace et al., 2019).	during manuscript preparation for language refine-	631
583	8 Ethical Considerations	ment, literature exploration, and \LaTeX formatting	632
584	Dual-use risk. This work studies six textual ob-	assistance. All generated content was carefully re-	633
585	fuscation strategies that could potentially be used	viewed, verified, and revised by the authors prior	634
586	to evade hate speech detection systems. However,	to submission.	635
587	these techniques are already widely documented	References	636
588	and commonly observed in real-world online dis-	Piush Aggarwal and Torsten Zesch. 2022. Analyzing	637
589	course (Gröndahl et al., 2018 ; Aggarwal and Zesch,	the real vulnerability of hate speech detection sys-	638
590	2022). Our contribution is not the creation of new	tems against targeted intentional noise . In <i>Proceed-</i>	639
591	attack mechanisms, but a systematic evaluation	<i>ings of the Eighth Workshop on Noisy User-generated</i>	640
592	framework for measuring model robustness against	<i>Text (W-NUT 2022)</i> , pages 230–242, Gyeongju, Re-	641
593	existing adversarial behaviour. The primary bene-	public of Korea. Association for Computational Lin-	642
594	ficiaries of this work are researchers and platform	guistics.	643
595	operators developing more reliable moderation sys-	Tommaso Caselli, Valerio Basile, Jelena Mitrović, and	644
596	tems.	Michael Granitzer. 2021. HateBERT: Retraining	645
597	Dataset release and reproducibility. To support	BERT for abusive language detection in English . In	646
598	reproducibility and future research in adversarial	<i>Proceedings of the 5th Workshop on Online Abuse</i>	647
599	hate speech detection, we plan to publicly release	<i>and Harms (WOAH 2021)</i> , pages 17–25, Online. As-	648
600	the obfuscation framework, generated dataset vari-	sociation for Computational Linguistics.	649
601	ants, and experimental code upon acceptance of the	Alexis Conneau, Kartikay Khandelwal, Naman Goyal,	650
602	paper. Although such resources could theoretically	Vishrav Chaudhary, Guillaume Wenzek, Francisco	651
603	be misused to automate toxic-text obfuscation, the	Guzmán, Edouard Grave, Myle Ott, Luke Zettle-	652
604	techniques studied are already publicly known and	moyer, and Veselin Stoyanov. 2020. Unsupervised	653
605	widely used online. We believe the scientific value	cross-lingual representation learning at scale . In <i>Pro-</i>	654
606	of transparent and standardised robustness evalua-	<i>ceedings of ACL</i> .	655
607	tion outweighs the associated risks, consistent with	Thomas Davidson, Dana Warmusley, Michael Macy, and	656
608	open-science principles adopted by the ACL com-	Ingmar Weber. 2017. Automated hate speech detec-	657
609	munity.	tion and the problem of offensive language . In <i>Pro-</i>	658
		<i>ceedings of the International AAAI Conference on</i>	659

660	<i>Web and Social Media (ICWSM)</i> , volume 11, pages	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint</i>	713
661	512–515.	<i>arXiv:2303.08774</i> .	714
662	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Paul Röttger, Bertie Vidgen, Dong Nguyen, and 1 others.	715
663	Kristina Toutanova. 2019. BERT: Pre-training of	2021. Hatecheck: Functional tests for hate speech	716
664	deep bidirectional transformers for language under-	detection models . In <i>Proceedings of ACL-IJCNLP</i> .	717
665	standing . In <i>Proceedings of NAACL</i> .		
666	Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing	Victor Sanh, Lysandre Debut, Julien Chaumond, and	718
667	Dou. 2018. HotFlip: White-box adversarial exam-	Thomas Wolf. 2019. Distilbert, a distilled version	719
668	ples for text classification . In <i>Proceedings of ACL</i> .	of bert: smaller, faster, cheaper and lighter . <i>ArXiv</i> ,	720
669	Jacob Eisenstein. 2013. What to do about bad language	abs/1910.01108 .	721
670	on the internet . In <i>Proceedings of NAACL</i> .		
671	Paula Fortuna and Sérgio Nunes. 2018. A survey of	Rico Sennrich, Barry Haddow, and Alexandra Birch.	722
672	automatic detection of hate speech in text . <i>ACM</i>	2016. Neural machine translation of rare words with	723
673	Computing Surveys , 51(4):1–30.	subword units . In <i>Proceedings of the 54th Annual</i>	724
674	Antigoni-Maria Founta, Constantinos Djouvas, De-	<i>Meeting of the Association for Computational Lin-</i>	725
675	spoina Chatzakou, Ilias Leontiadis, Jeremy Black-	<i>guistics (ACL)</i> , pages 1715–1725.	726
676	burn, Gianluca Stringhini, Athena Vakali, Michael	Sali A. Tagliamonte and Derek Denis. 2008. Linguistic	727
677	Sirivianos, and Nicolas Kourtellis. 2018. A large	ruin? LOL! instant messaging and teen language .	728
678	scale crowdsourcing dataset for abusive language de-	<i>American Speech</i> , 83(1).	729
679	tection . In <i>Proceedings of ICWSM</i> .		
680	Tommi Gröndahl, Luca Pajola, Mika Juuti, Mauro Conti,	Hugo Touvron, Thibaut Lavril, Gautier Izacard, and 1	730
681	and N. Asokan. 2018. All you need is "love": Evad-	others. 2023. Llama: Open and efficient foundation	731
682	ing hate speech detection . In <i>Proceedings of the</i>	language models . <i>arXiv preprint arXiv:2302.13971</i> .	732
683	<i>11th ACM Workshop on Artificial Intelligence and</i>	Bertie Vidgen and Leon Derczynski. 2020. Direc-	733
684	<i>Security</i> .	tions in abusive language training data, a system-	734
685	Bo Han and Timothy Baldwin. 2011. Lexical normali-	atic review: Garbage in, garbage out . <i>PLOS ONE</i> ,	735
686	sation of short text messages: Makn sens a #twitter .	15(12):e0243300.	736
687	In <i>Proceedings of ACL</i> .	Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gard-	737
688	Laura Hanu and Unitary Team. 2020. Detoxify .	ner, and Sameer Singh. 2019. Universal adversarial	738
689	Ipsos. 2023. Two in three people often encounter hate	triggers for attacking and analyzing NLP . In <i>Pro-</i>	739
690	speech online .	<i>ceedings of EMNLP</i> .	740
691	Erik Jones, Robin Jia, Aditi Raghunathan, and Percy	Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan,	741
692	Liang. 2020. Robust encodings: A framework for	Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah,	742
693	combating adversarial typos . In <i>Proceedings of ACL</i> .	and Bo Li. 2021. AdvGLUE: A multi-task bench-	743
694	Hannah Kirk, Bertie Vidgen, Paul Röttger, and 1 oth-	mark for robustness evaluation of language models .	744
695	ers. 2022. Hatemoji: A test suite and adversarially-	In <i>Proceedings of NeurIPS</i> .	745
696	generated dataset for benchmarking and detecting	Zeerak Waseem and Dirk Hovy. 2016. Hateful symbols	746
697	emoji-based hate . In <i>Proceedings of NAACL-HLT</i> .	or hateful people? predictive features for hate speech	747
698	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	detection on Twitter . In <i>Proceedings of the NAACL</i>	748
699	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	<i>Student Research Workshop</i> , pages 88–93, San Diego,	749
700	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	California. Association for Computational Linguis-	750
701	Roberta: A robustly optimized bert pretraining ap-	<i>tics</i> .	751
702	proach . <i>ArXiv</i> , abs/1907.11692.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	752
703	Meta Platforms. 2024. Facebook hate speech removal	Chaumond, Clément Delangue, Anthony Moi, Pierric	753
704	statistics .	Cistac, and 1 others. 2020. Transformers: State-of-	754
705	Sandip Modha, Prasenjit Majumder, and Thomas Mandl.	the-art natural language processing . In <i>Proceedings</i>	755
706	2020. Detecting offensive speech in conversational	<i>of EMNLP: System Demonstrations</i> , pages 38–45.	756
707	code-switching . In <i>Proceedings of LREC</i> .		
708	John Morris, Eli Lifland, Jin Yong Yoo, and 1 others.	A Appendix	757
709	2020. Textattack: A framework for adversarial at-	A.1 Hyperparameter Settings	758
710	tacks, data augmentation, and adversarial training in	All experiments were conducted using pretrained	759
711	nlp . In <i>Proceedings of EMNLP: System Demonstra-</i>	transformer-based models implemented with the	760
712	<i>tions</i> .	HuggingFace Transformers framework and Py-	761
		Torch on Google Colab (NVIDIA T4 GPU). The	762
		models were evaluated using a learning rate of $2 \times$	763
		10^{-5} , batch size 16, and maximum sequence length	764

765 of 128 tokens. Consistent tokenization, padding,
766 and truncation settings were applied across all ex-
767 periments to ensure fair comparison between clean
768 and obfuscated hate speech datasets. F1 Score was
769 used as the primary evaluation metric.