

Perturbations in the Wild: Leveraging Human-Written Text Perturbations for Realistic Adversarial Attack and Defense

Anonymous ACL submission

Abstract

We propose a novel algorithm, ANTHRO, that *inductively* extracts over 600K human-written text perturbations in the wild and leverages them for *realistic* adversarial attack. Unlike existing character-based attacks which often deductively hypothesize a set of manipulation strategies, our work is grounded on actual observations from real-world texts. We find that adversarial texts generated by ANTHRO achieve the best trade-off between (1) attack success rate, (2) semantic preservation of the original text, and (3) stealthiness—i.e. indistinguishable from human writings hence harder to be flagged as suspicious. Specifically, our attacks accomplished around 83% and 91% attack success rates on BERT and RoBERTa, respectively. Moreover, it outperformed the *TextBugger* baseline with an increase of 50% and 40% in terms of semantic preservation and stealthiness when evaluated by both layperson and professional human workers. ANTHRO can further enhance a BERT classifier’s performance in understanding different variations of human-written toxic texts via adversarial training when compared to the Perspective API. *All source code will be released.*

1 Introduction

Machine learning (ML) models trained to optimize only the prediction performance are often vulnerable to adversarial attacks (Papernot et al., 2016; Wang et al., 2019). In the text domain, especially, a character-based adversarial attacker aims to fool a target ML model by generating an adversarial text x^* from an original text x by manipulating characters of different words in x , such that some properties of x are preserved (Li et al., 2018; Eger et al., 2019; Gao et al., 2018). We characterize strong and practical adversarial attacks as three criteria: (1) *attack performance*, as measured by the ability to flip a target model’s predictions, (2) *semantic preservation*, as measured by the ability

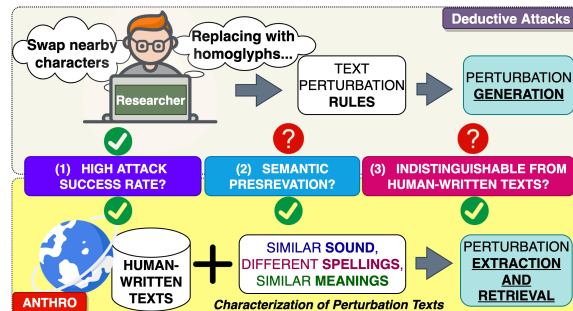


Figure 1: ANTHRO (Bottom) extracts and uses human-written perturbations for adversarial attacks instead of proposing a specific set of manipulation rules (Top).

to preserve the meaning of an original text, and (3) *stealthiness*, as measured by how unlikely it is detected as machine-manipulation and removed by defense systems or human examiners (Figure 1). While the first two criteria are natural derivation from adversarial literature (Papernot et al., 2016), stealthiness is also important to be a practical attack under a mass-manipulation scenario.

Previously proposed character-based attacks follow a *deductive* approach where the researchers hypothesize a set of text manipulation strategies that exploit some vulnerabilities of textual ML models (Figure 1). Although these deductively derived techniques can demonstrate superior attack performance, there is no guarantee that they also perform well with regard to semantic preservation and stealthiness. We first analyze why enforcing these properties are challenging especially for character-based attacks.

To preserve the semantic meanings, an attacker can minimize the distance between representative vectors learned from a large pre-trained model—e.g., Universal Sentence Encoder (Cer et al., 2018) of the two sentences. However, this is only applicable in word- or sentence-based attacks, not in character-based attacks. It is because character-based manipulated tokens are more prone to become out-of-distribution—e.g., morons→mor0ns, from what is observed in a typical training cor-

Attacker #texts, #tokens	Reddit Comts. »5B, N/A	News Comts. (34M, 11M)
TextBugger	51.6% (126/244)	7.10% (11K/152K)
VIPER	3.2% (1/31)	0.13% (25/19K)
DeepWordBug	0% (0/31)	0.27% (51/19K)
ANTHRO	82.4% (266/323)	55.7% (16K/29K)

Table 1: Percentage of offensive perturbed words generated by different attacks that can be observed in real human-written comments on Reddit and online news.

personal exposure from Reddit or YouTube comments to decide if a word choice looks natural (Sec. 4.2). Quantitatively, we discover that not all the perturbations generated by deductive methods are observed on the Web (Table 1). To analyze this, we first use each attack to generate all possible perturbations of either (1) a list of over 3K unique offensive words or (2) a set of the top 5 offensive words (“c*nt”, “b*tch”, “m*therf***er”, “bast*rd”, “d*ck”). Then, we calculate how many of the perturbed words are present in a dataset of over 34M online news comments or are used by at least 50 unique commentators on Reddit, respectively. Even though *TextBugger* was well-known to simulate human-written typos as adversarial texts, merely 51.6% and 7.1% of its perturbations are observed on Reddit and online news comments, implying *TextBugger*’s generated adversarial texts being “unnatural” and “easily-detectable” by human-in-the-loop defense systems.

2.2 The SMS Property: Similar Sound, Similar Meaning, Different Spelling

The existence of a non-arbitrary relationship between sounds and meanings has been proven by a life-long research establishment (Köhler, 1967; Jared and Seidenberg, 1991; Gough et al., 1972). In fact, Blasi et al. (2016) analyzed over 6K languages and discovered a high correlation between a word’s sound and meaning both inter- and intra-cultures. Aryani et al. (2020) found that how a word sounds links to an individual’s emotion. This motivates us to hypothesize that words spelled differently yet have the same meanings such as text perturbations will also have similar sounds.

Figure 2 displays several perturbations that are found from real-life texts. Even though these perturbations are *spelled differently* from the original word, they all preserve *similar meanings* when perceived by humans. Such semantic preservation is feasible because humans perceive these variations *phonetically similar* to the respective origi-

nal words (Van Orden, 1987). For example, both “republican” and “republikan” sound similar when read by humans. Therefore, given the surrounding context of a perturbed sentence—e.g., “*President Trump is a republican*”, and the phonetic similarity of “republican” and “republikan”, end-users are more likely to interpret the perturbed sentence as “*President Trump is a republican*”. We call these characteristics of text perturbations the *SMS* property: “*similar Sound, similar Meaning, different Spellings*”. Noticeably, the *SMS* characterization includes a subset of “visually similar” property of perturbations as studied in previous adversarial attacks such as *TextBugger* (e.g., “hello” sounds similar with “hello”), *VIPER* and *DeepWordBug*. However, two words that look very similar sometimes carry different meanings—e.g., “garbage”→“gabrage”. Moreover, our characterization is also distinguished from *homophones* (e.g., “to” and “two”) which describe words with similar sound yet *different meaning*.

3 A Realistic Adversarial Attack

Given the above analysis, we now derive our proposed ANTHRO adversarial attack. We first share how to systematically encode the sound—i.e., phonetic feature, of any given words and use it to search for their human-written perturbations that satisfy the *SMS* property. Then, we introduce an iterative algorithm that utilizes the extracted perturbations to attack textual ML models.

3.1 Mining Perturbations in the Wild

Sound Encoding with SOUNDEX++. To capture the sound of a word, we adopt and extend the case-insensitive SOUNDEX algorithm. SOUNDEX helps index a word based on how it sounds rather than how it is spelled (Stephenson, 1980). Given a word, SOUNDEX first keeps the 1st character. Then, it removes all vowels and matches the remaining characters *one by one* to a digit following a set of predefined rules—e.g., “B”, “F”→1, “D”, “T”→3 (Stephenson, 1980). For example, “Smith” and “Smyth” are both encoded as S530.

As the SOUNDEX system was designed mainly for encoding surnames, it does not necessarily work for texts in the wild. For example, it cannot recognize visually-similar perturbations such as “l”→“1”, “a”→“@” and “O”→“0”. Moreover, it always fixes the 1st character as part of the final encodes. This rule is too rigid and can result in words

Word	SOUNDEX	SOUNDEX++ (Ours)
porn	P650	P650 (k=0), PO650 (k=1)
p0rn	P065(X)	(same as above)
lesbian	L215	L245 (k=0), LE245 (k=1)
lesbbi@n	L21@(X)	(same as above)
losbian	L215(X)	L245 (k=0), LO245 (k=1)

(X): Incorrect encoding

Table 2: SOUNDEX++ can capture visually similar characters and is more accurate in differentiating between desired (blue) and undesired (red) perturbations.

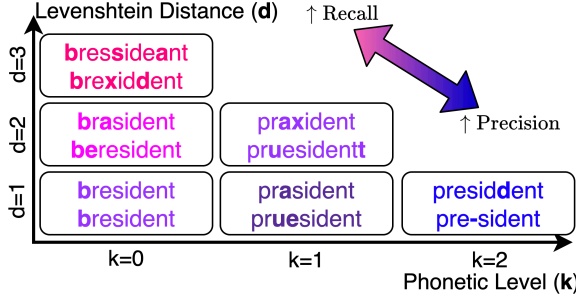


Figure 3: Trade-off between precision and recall of extracted perturbations for the word “president” w.r.t different k and d values. Higher k and lower d associate with better preservation of the original meaning.

that are entirely different yet encoded the same (Table 2). To solve these issues, we propose a new SOUNDEX++ algorithm. Not only SOUNDEX++ encodes visually-similar characters the same, it also encodes the sound of a word at different hierarchical levels k (Table 2). At level $k=0$, SOUNDEX++ works similar to SOUNDEX by fixing the first character. At level $k \geq 1$, SOUNDEX++ instead fixes the first $k+1$ characters and encodes the rest.

Levenshtein Distance d and Phonetic Level k as a Semantic Preservation Proxy. Since SOUNDEX++ is not designed to capture a word’s semantic meaning, we utilize both phonetic parameter k and *Levenshtein distance* d (Levenshtein et al., 1966) as a heuristic approximation to measure the semantic preservation between two words. Intuitively, the higher the phonetic level ($k \geq 1$) at which two words share the same SOUNDEX++ code and the smaller the Levenshtein distance d to transform one word to another, the more likely human associate them with the meaning. In other words, k and d are hyper-parameters that help control the trade-off between precision and recall when retrieving perturbations of a given word such that they satisfy the SMS property (Figure 3). We will later carry out a human study to evaluate how well our extracted perturbations can preserve the

Algorithm 1 ANTHRO Attack Algorithm

```

1: Input:  $\{H\}_0^K, k, d$ 
2: Input: target classifier  $f$ , original sentence  $x$ 
3: Output: perturbed sentence  $x^*$ 
4: Initialize:  $x^* \leftarrow x$ 
5: for word  $x_i$  in  $x$  do:  $s_i \leftarrow \text{Score}(x_i, f)$ 
6:  $\mathcal{W}_{\text{order}} \leftarrow \text{Sort}(x_1, x_2, \dots, x_m)$  according to  $s_i$ 
7: for  $x_i$  in  $\mathcal{W}_{\text{order}}$  do:
8:    $\mathcal{P} \leftarrow \text{ANTHRO}(x_i, k, d, \{H\}_0^K)$  // Eq.(3)
9:    $x^* \leftarrow$  replace  $x_i \in x$  with the best  $w \in \mathcal{P}$ 
10:  if  $f(x^*) \neq f(x)$  then return  $x^*$ 
11: return None

```

semantic meanings in practice.

Mining from the Wild. To mine all human-written perturbations, we first collect a large corpus \mathcal{D} of over 18M sentences written by netizens from 9 different datasets (Table A.1 in Appendix). We select these datasets because they include offensive texts such as hate speech, sensitive search queries, etc., and hence very likely to include text perturbations. Next, for each phonetic level $k \leq K$, we curate different hash tables $\{H\}_0^K$ that maps a unique SOUNDEX++ code c to a set of its matching unique *case-sensitive* tokens that share the same encoding c as follows:

$$H_k : c \mapsto \{w_j | S(w_i, k) = S(w_j, k) = c \quad \forall w_i, w_j \in \mathcal{D}, w_i \neq w_j\}, \quad (1)$$

where $S(w, k)$ returns the SOUNDEX++ code of token w at phonetic level k , K is the largest phonetic level we want to encode. With $\{H\}_0^K, k$ and d , we can now search for the set of perturbations $G_k^d(w^*)$ of a specific target token w^* as follows:

$$G_k^d(w^*) \leftarrow \{w_j | w_j \in H_k[S(w^*, k)], \text{Lev}(w^*, w_j) \leq d\} \quad (2)$$

where $\text{Lev}(w^*, w_j)$ returns the Levenshtein distance between w^* and w_j . Noticeably, we only extract $\{H\}_0^K$ **once** from \mathcal{D} via Eq. (1), then we can use Eq. (2) to retrieve all perturbations for a given word during deployment. We name this method of mining and retrieving human-written text perturbations in the wild as **ANTHRO**, aka *human-like* perturbations:

$$\text{ANTHRO} : w^*, k, d, \{H\}_0^K \mapsto G_k^d(w^*) \quad (3)$$

ANTHRO Attack. To utilize ANTHRO for adversarial attack on model $f(x)$, we propose the ANTHRO attack algorithm (Alg. 1). We use the

same iterative mechanism (Ln.9–13) that is common among other black-box attacks. This process replaces the most vulnerable word in sentence x , which is evaluated with the support of $\text{Score}(\cdot)$ function (Ln. 5), with the perturbation that best drops the prediction probability $f(x)$ on the correct label. Unlike the other methods, ANTHRO inclusively draws from perturbations extracted from human-written texts captured in $\{\mathcal{H}\}_0^K$ (Ln. 10). We adopt the $\text{Score}(\cdot)$ from *TextBugger*.

4 Evaluation

We evaluate ANTHRO by: (1) attack performance, (2) semantic preservation, and (3) Turing Test (TT)—i.e., how likely an attack message is spotted as machine-generated by human examiners.

4.1 Attack Performance

Setup. We use BERT (*case-insensitive*) (Jin et al., 2019) and RoBERTa (*case-sensitive*) (Liu et al., 2019) as target classifiers to attack. We evaluate on three public tasks, namely detecting toxic comments ((TC) dataset, Kaggle 2018), hate speech ((HS) dataset (Davidson et al.)), and online cyberbullying texts ((CB) dataset (Wulczyn et al., 2017a)). We split each dataset to *train*, *validation* and *test* set with the 8:1:1 ratio. Then, we use the train set to fine-tune BERT and RoBERTa with a maximum of 3 epochs and select the best checkpoint using the validation set. BERT and RoBERTa achieve around 0.85–0.97 in F1 score on the test sets (Table A.2 in Appendix). We evaluate with targeted attack (change positive→negative label) since it is more practical. We randomly sample 200 examples from each test set and use them as initial sentences to attack. We repeat the process 3 times with unique random seeds and report the results. We use the *attack success rate* (Atk%) metric—i.e., the number of examples whose labels are flipped by an attacker over the total number of texts that are correctly predicted pre-attack. We use the 3rd party open-source *OpenAttack* (Zeng et al., 2021) framework to run all evaluations.

Baselines. We compare ANTHRO with three baselines, namely *TextBugger* (Li et al., 2018), *VIPER* (Eger et al., 2019) and *DeepWordBug* (Gao et al., 2018). These attackers utilize different character-based manipulations to craft their adversarial texts as described in Sec. 1. From the analysis in Sec. 3.1 and Figure 3, we set $k \leftarrow 1$ and

$d \leftarrow 1$ for ANTHRO to achieve a balanced trade-off between precision and recall on the SMS property. We examine all attackers under several combinations of different normalization layers. They are (1) *Accents normalization* (A) and (2) *Homoglyph normalization*¹ (H), which converts non-English accents and homoglyphs to their corresponding ascii characters, (3) *Perturbation normalization* (P), which normalizes potential character-based perturbations using the SOTA misspelling correction model *Neuspell* (Jayanthi et al., 2020). These normalizers are selected as counteracts against the perturbation strategies employed by *VIPER* (uses non-English accents), *DeepWordBug* (uses homoglyphs) and *TextBugger*, ANTHRO (based on misspelling and typos), respectively.

Results. Overall, both ANTHRO and *TextBugger* perform the best, with ANTHRO being the most robust attacker on RoBERTa (due to its case-sensitive perturbations) and is competitive compared to *TextBugger* on BERT (Table 3). Because RoBERTa uses the accent \ddot{G} as a part of its byte-level BPE encoding (Liu et al., 2019) while BERT by default removes all non-English accents, *VIPER* achieves a near perfect score on RoBERTa, yet it is ineffective on BERT. Since *VIPER* exclusively utilizes non-English accents, its attacks can be easily corrected by the *accents normalizer* (Table 3). Similarly, *DeepWordBug* perturbs texts with homoglyph characters, most of which can also be normalized using a 3rd party homoglyph detector (Table 3).

In contrast, even under all normalizers—i.e., A+H+P, *TextBugger* and ANTHRO still achieves 66.3% and 73.7% in Atk% on average across all evaluations. Although *Neuspell* (Jayanthi et al., 2020) drops *TextBugger*’s Atk% 14.7% across all runs, it can only reduce the Atk% of ANTHRO a mere 7.5% on average. This is because *TextBugger* and *Neuspell* or other dictionary-based typo correctors rely on fixed deductive rules—e.g., swapped, replaced by neighbor letters, for attack and defense. However, ANTHRO utilizes human-written perturbations which are greatly varied, hence less likely to be systematically detected.

4.2 Semantic Preservation and Turing Test

Since ANTHRO and *TextBugger* are the top two effective attacks, this section will focus on evaluating their ability in semantic preservation and

¹ <https://github.com/codebox/homoglyph>

Attacker	Normalizer	BERT (<i>case-insensitive</i>)			RoBERTa (<i>case-sensitive</i>)		
		Toxic Comments	HateSpeech	Cyberbullying	Toxic Comments	HateSpeech	Cyberbullying
TextBugger	-	0.76±0.02	0.94±0.01	0.78±0.03	0.77±0.06	0.87±0.01	0.72±0.01
DeepWordBug	-	0.56±0.04	0.68±0.01	0.50±0.02	0.52±0.01	0.42±0.04	0.38±0.04
VIPER	-	0.08±0.03	0.01±0.01	0.13±0.02	1.00±0.00	1.00±0.00	0.99±0.01
ANTHRO	-	0.72±0.02	0.82±0.01	0.71±0.02	0.84±0.00	0.93±0.01	0.78±0.01
TextBugger	A	-	-	-	0.72±0.02	0.92±0.00	0.74±0.02
DeepWordBug	A	-	-	-	0.43±0.02	0.59±0.03	0.43±0.01
VIPER	A	-	-	-	0.09±0.01	0.05±0.01	0.17±0.02
ANTHRO	A	-	-	-	0.77±0.02	0.94±0.02	0.84±0.02
TextBugger	A+H	0.78±0.03	0.85±0.00	0.79±0.00	0.74±0.02	0.93±0.01	0.77±0.03
DeepWordBug	A+H	0.04±0.00	0.06±0.02	0.01±0.01	0.03±0.01	0.01±0.01	0.06±0.02
VIPER	A+H	0.07±0.00	0.01±0.01	0.10±0.00	0.13±0.02	0.07±0.01	0.17±0.01
ANTHRO	A+H	0.76±0.02	0.77±0.03	0.73±0.05	0.82±0.02	0.97±0.00	0.82±0.02
TextBugger	A+H+P	0.73±0.02	0.64±0.06	0.70±0.04	0.68±0.06	0.57±0.03	0.66±0.04
DeepWordBug	A+H+P	0.02±0.01	0.04±0.02	0.01±0.01	0.02±0.01	0.01±0.01	0.02±0.01
VIPER	A+H+P	0.12±0.01	0.04±0.01	0.17±0.03	0.11±0.02	0.05±0.01	0.18±0.01
ANTHRO	A+H+P	0.65±0.04	0.64±0.01	0.60±0.05	0.80±0.02	0.91±0.03	0.82±0.02

(-) BERT already has the accents normalization (A normalizer) by default, (Red): Poor performance (Atk%<0.15)

Table 3: Averaged attack success rate (Atk%↑) of different attack methods

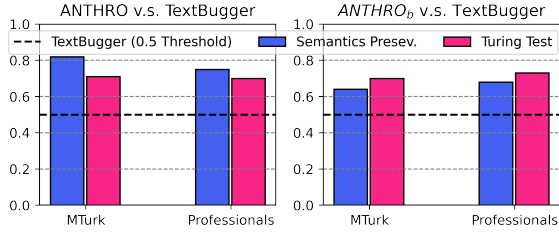


Figure 4: Semantic preservation and Turing test results

Turing test. Given an original sentence x and its adversarial text x^* generated by either one of the attacks, we design a human study to directly compare ANTHRO with *TextBugger*. Specifically, two alternative hypotheses for our validation are (1) $\mathcal{H}_{\text{Semantics}}$: x^* generated by ANTHRO preserves the original meanings of x *better* than that generated by *TextBugger* and (2) $\mathcal{H}_{\text{Turing}}$: x^* generated by ANTHRO is *more likely* to be perceived as a human-written text (and not machine) than that generated by *TextBugger*, hence called Turing test (TT) (Uchendu et al., 2020).

Human Study Design. We use the two attackers to generate adversarial texts targeting BERT model on 200 examples sampled from the TC dataset’s test set. We then gather examples that are successfully attacked by both ANTHRO and *TextBugger*. Next, we present a pair of texts, one generated by ANTHRO and one by *TextBugger*, together with the original sentence to human subjects. We then ask them to select (1) which text better preserves the meaning of the original sentence (Figure A.2 in Appendix) and (2) which text is more likely to be written by human (Figure A.3

Reason	Favorable From ANTHRO	Unfavorable From TextBugger
Genuine Typos	stuupid, but, Faoggt	sutpid, burt, Foggat
Intelligible	failiure	faioure
Sound Preserv.	shytty, crp	shtty, crsp
Meaning Preserv.	ga-y, ashole, dummb	bay, alshose, dub
High Search Results	sodmized, kiills	Smdooized, klils
Personal Exposure	ign0rant, gaarbage	ignorajt, garage
Word Selection	morons→mor0ns	edited→ewited

Table 4: Top reasons in favoring ANTHRO’s perturbations as more likely to be written by human.

in Appendix). To reduce bias, we present the two questions in two separate tasks. The human subjects include both (1) Amazon Mechanical Turk (MTurk) workers and (2) professional data annotators at a company with extended experience in annotating texts in domain such as toxic and hate speech. Our human subject study with MTurk workers was IRB-approved.

Quantitative Results. It is statistically significant ($p\text{-value} \leq 0.05$) to reject the null hypotheses of both $\mathcal{H}_{\text{Semantics}}$ and $\mathcal{H}_{\text{Turing}}$ (Table A.3 in Appendix). Overall, adversarial texts generated by perturbations mined in the wild are much better at preserving the original semantics and also more indistinguishable from human-written texts than those generated by *TextBugger* (Figure 4, Left).

Qualitative Analysis. We also ask the professional subjects to provide optional comments on their thought process. Table 4 summarizes the top reasons why they favor ANTHRO over *TextBugger* in terms of Turing test. ANTHRO’s perturbations

Attacker	Normalizer	BERT (case-insensitive)			RoBERTa (case-sensitive)		
		Toxic Comments	HateSpeech	Cyberbullying	Toxic Comments	HateSpeech	Cyberbullying
TextBugger	-	0.76±0.02	0.94±0.01	0.78±0.03	0.77±0.06	0.87±0.01	0.72±0.01
ANTHRO _β	-	0.82±0.01	0.97±0.01	0.88±0.04	0.91±0.02	0.97±0.01	0.89±0.02
TextBugger	A+H+P	0.73±0.02	0.64±0.06	0.70±0.04	0.68±0.06	0.57±0.03	0.66±0.04
ANTHRO _β	A+H+P	0.85±0.04	0.79±0.02	0.84±0.03	0.88±0.04	0.93±0.01	0.91±0.01

Table 5: Averaged attack success rate (Atk%↑) of ANTHRO_β and *TextBugger*

are perceived similar to genuine typos and more intelligible (than ones that might be generated by machine). They also better preserve both meanings and sounds. Moreover, some annotators also rely on personal exposure on Reddit, YouTube comments, or the frequency of word use via the search function on Reddit to decide if a word-choice is human-written. Interestingly, one mentions that ANTHRO is better at selecting sensible words—i.e., “morons” instead of “edit”, to perturb than *TextBugger*, even though the two methods share the same iterative attack mechanism (Alg. 1). This happens because ANTHRO directly ensembles the distribution of human-written texts, which naturally includes more replacement candidates for offensive than non-offensive words. This eventually increases the probability of sensitive words being perturbed.

5 ANTHRO_β Attack

ANTHRO_β. We want to examine if perturbations inductively extracted from the wild can help augment a deductive attack such as *TextBugger* and improve its overall performance. Hence, we introduce ANTHRO_β, which considers the perturbation candidates from both ANTHRO and *TextBugger* in Ln. 10 of Alg. 1. Alg. 1 still selects the perturbation that best flip the target model’s prediction.

Attack Performance. Even though ANTHRO comes second after *TextBugger* when attacking BERT model, Table 5 shows that when combined with *TextBugger*—i.e., ANTHRO_β, it consistently achieves superior performance with an average of 82.7% and 90.7% in Atk% on BERT and RoBERTa even under all normalizers (A+H+P).

Semantic Preservation and TT. ANTHRO_β improves *TextBugger*’s Atk% to over 8% on average (Table 5). It also improves *TextBugger*’s semantic preservation and TT score 32% and 42% (from 0.5 threshold) (Figure 4, Right). The presence of only a few human-like perturbations generated by ANTHRO is sufficient to signal whether

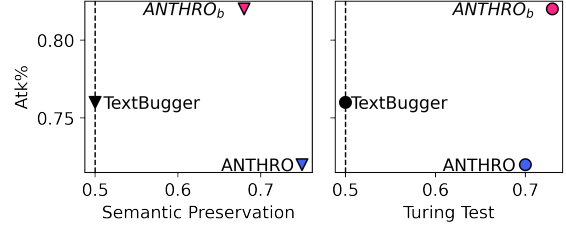


Figure 5: Trade-off among evaluation metrics

or not it is written by humans, while only one unreasonable perturbation generated by *TextBugger* can adversely affect its meaning. This explains the performance drop in terms of semantic preservation but not in TT when indirectly comparing ANTHRO_β with ANTHRO. Overall, ANTHRO_β also has the best trade-off between Atk% and human evaluation—i.e., positioning at top right corners in Figure 5. Particularly, ANTHRO_β trade-offs from ANTHRO some reduction in semantic preservation for superior Atk%. This gain iterates the overall benefits of human-written perturbations for adversarial attacks.

6 Defend ANTHRO, ANTHRO_β Attack

We examine two approaches to defend against ANTHRO attack. We compare them against BERT and BERT combined with 3 layers of normalization A+H+P. BERT is selected as it is better than RoBERTa at defending against ANTHRO (Table 3).

Sound-Invariant Textual Model (SOUNDCNN):

When the defender do *not* have access to the hash tables $\{\mathcal{H}\}_0^K$ of the attacker, the defender can train a generic model that encodes not the spellings but the phonetic features of a text for prediction. As an example, we train a CNN model (Kim, 2014) on top of a continuous embeddings layer for discrete SOUNDEX++ encodings of each token in a sentence. **Adversarial Training (ADV.TRAIN):** To overcome the lack of access to $\{\mathcal{H}\}_0^K$, the defender can extract his/her perturbations in the wild from a separate corpus \mathcal{D}^* where $\mathcal{D}^* \cap \mathcal{D} = \emptyset$ and use them to augment the training examples—i.e., via self-attack with ratio 1:1, to fine-tune a more

Model	ANTHRO			ANTHRO _{β}		
	TC \downarrow	HS \downarrow	CB \downarrow	TC \downarrow	HS \downarrow	CB \downarrow
BERT	0.72	0.82	0.71	0.82	0.97	0.88
BERT+A+H+P	0.65	0.65	0.60	0.85	0.79	0.84
ADV.TRAIN	0.41	0.30	0.35	0.72	0.72	0.67
SOUNDCNN	0.14	0.02	0.15	0.86	0.84	0.92

Table 6: Averaged Atk% \downarrow of ANTHRO and ANTHRO _{β} against different defense models.

robust BERT model. Here we use \mathcal{D}^* as a corpus of 34M general comments from online news.

Results. We follow the same evaluation procedure in Sec. 4.1. Table 6 shows that both SOUNDCNN and ADV.TRAIN are robust against ANTHRO attack, while ADV.TRAIN performs best when defending ANTHRO _{β} . Since SOUNDCNN is strictly based on phonetic features, it is vulnerable against ANTHRO _{β} whenever *TextBugger*’s perturbations are selected. Table 6 also underscores that ANTHRO _{β} is a strong and practical attack, defense against which is thus an important future direction.

7 Discussion

Evaluation with *Perspective API*. The understanding of different variations of human-written texts is critical to fully capture the semantic meanings of inputs especially in sensitive domains such as toxicity moderation. We utilize ANTHRO and ANTHRO _{β} to evaluate such robust understanding of the popular *Perspective API*², which has been adopted in various publishers—e.g., NYTimes, and platforms—e.g., Disqus, Reddit. Specifically, we evaluate (1) if the API can capture different forms of human-written toxic texts by using ANTHRO to randomly perturb different portions of words in 200 positive texts from the TC dataset, (2) if the API can defend against ANTHRO attack either via a direct iteration mechanism (Alg. 1) or via transfer attack through an intermediate BERT classifier.

Figure 6 (Left) shows that the API service provides superior performance compared to a self fine-tuned BERT classifier, yet its precision deteriorates quickly from 0.95 to only 0.9 and 0.82 when 25%–50% of a sentence are perturbed. In contrast, the ADV.TRAIN (Sec. 7) model achieves fairly consistent precision in the same setting. The API is also more vulnerable against both direct and transfer attacks from our proposed attacks (Figure. 6, Right) than *TextBugger*, with its preci-

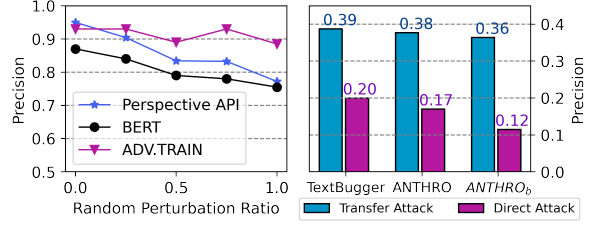


Figure 6: (Left) Precision on human-written perturbed texts synthesized by ANTHRO and (Right) Robustness evaluation of *Perspective API* under different attacks

sion dropped to only 0.12 when evaluated against ANTHRO _{β} . Overall, not only ANTHRO is a powerful and realistic attack, it can also help develop more robust text classifiers in practice.

Computational Complexity. The **one-time** extraction of $\{\mathcal{H}\}_0^K$ via Eq. (1) has $\mathcal{O}(|\mathcal{D}|L)$ where $|\mathcal{D}|$, L is the # of tokens and the length of longest token in \mathcal{D} (hash-map operations cost $\mathcal{O}(1)$). Given a word w and \mathbf{k}, \mathbf{d} , ANTHRO retrieves a list of perturbation candidates via Eq. (2) with $\mathcal{O}(|w|\max(\mathcal{H}_k))$ where $|w|$ is the length of w and $\max(\mathcal{H}_k)$ is the size of the largest set of tokens sharing the same SOUNDEX++ encoding in \mathcal{H}_k . Since $\max(\mathcal{H}_k)$ is constant, the upper-bound then becomes $\mathcal{O}(|w|)$.

Limitation The perturbation candidate retrieval operation (Eq. (2)) has a higher computational complexity than that of other methods—i.e., $\mathcal{O}(|w|)$ v.s. $\mathcal{O}(1)$. This can prolong the running time, especially when attacking long documents. However, we can overcome this by storing all the perturbations (given \mathbf{k}, \mathbf{d}) of the top frequently used offensive and non-offensive English words. We can then expect the operation to have an average complexity close to $\mathcal{O}(1)$. The current SOUNDEX++ algorithm is designed for English texts and might not be applicable in other languages. Thus, we plan to extend ANTHRO to a multilingual setting.

8 Conclusion

We propose ANTHRO, a character-based attack algorithm that extracts human-written perturbations in the wild and then utilizes them for adversarial text generation. Our approach yields the best trade-off between attack performance, semantic preservation and stealthiness under both empirical experiments and human studies. A BERT classifier trained with examples augmented by ANTHRO can also better understand human-written texts.

² <https://www.perspectiveapi.com/>

References

- Arash Aryani, Erin S Isbilen, and Morten H Christiansen. 2020. Affective arousal links sound to meaning. *Psychological science*, 31(8):978–986.
- Damián E Blasi, Søren Wichmann, Harald Hammarström, Peter F Stadler, and Morten H Christiansen. 2016. Sound–meaning association biases evidenced across thousands of languages. *Proceedings of the National Academy of Sciences*, 113(39):10818–10823.
- Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *EMNLP’18, Demo*.
- Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. Automated hate speech detection and the problem of offensive language. In *ICWSM’17*.
- Steffen Eger, Gözde Gül Şahin, Andreas Rücklé, Ji-Ung Lee, Claudia Schulz, Mohsen Mesgar, Krishnakant Swarnkar, Edwin Simpson, and Iryna Gurevych. 2019. [Text processing like humans do: Visually attacking and shielding NLP systems](#). In *NAACL’19*, pages 1634–1647, Minneapolis, Minnesota.
- Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun Qi. 2018. Black-box generation of adversarial text sequences to evade deep learning classifiers. In *SPW’18*. IEEE.
- Raul Gomez, Jaume Gibert, Lluís Gomez, and Dimosthenis Karatzas. 2020. Exploring hate speech detection in multimodal publications. In *WACV’20*, pages 1470–1478.
- Philip B Gough, JF Kavanagh, and IG Mattingly. 1972. One second of reading. *Cambridge: MIT Press*, pages 331–358.
- Debra Jared and Mark S Seidenberg. 1991. Does word identification proceed from spelling to sound to meaning? *Journal of Experimental Psychology: General*, 120(4):358.
- Sai Muralidhar Jayanthi, Danish Pruthi, and Graham Neubig. 2020. NeuSpell: A neural spelling correction toolkit. In *EMNLP’20, Demo*, Online.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2019. Is bert really robust? natural language attack on text classification and entailment. *arXiv preprint arXiv:1907.11932*.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *EMNLP’14*, Doha, Qatar.
- Elena Kochkina, Maria Liakata, and Arkaitz Zubiaga. 2018. [All-in-one: Multi-task learning for rumour verification](#). In *ACL’18*, Santa Fe, New Mexico, USA. ACL.
- Wolfgang Köhler. 1967. Gestalt psychology. *Psychologische Forschung*, 31(1):XVIII–XXX.
- Thai Le, Suhan Wang, and Dongwon Lee. 2020. Malcom: Generating malicious comments to attack neural fake news detection models. In *ICDM’20*. IEEE.
- Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union.
- Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2018. TextBugger: Generating Adversarial Text Against Real-world Applications. *NDSS’18*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. 2016. The limitations of deep learning in adversarial settings. In *EuroS&P’16*, pages 372–387. IEEE.
- Danish Pruthi, Bhuwan Dhingra, and Zachary C Lipton. 2019. Combating adversarial misspellings with robust word recognition. In *ACL’19*.
- Charles Stephenson. 1980. The methodology of historical census record linkage: A user’s guide to the soundex. *Journal of Family History*, 5(1):112–115.
- Caroline Tagg. 2011. Wot did he say 01" could u not c him 4 dust?: Written and spoken creativity in text messaging. *Transforming literacies and language: Multimodality and literacy in the new media age*, 223.
- Adaku Uchendu, Thai Le, Kai Shu, and Dongwon Lee. 2020. [Authorship attribution for neural text generation](#). In *EMNLP’20*, pages 8384–8395, Online. Association for Computational Linguistics.
- Guy C Van Orden. 1987. A rows is a rose: Spelling, sound, and reading. *Memory & cognition*, 15(3):181–198.
- Wenqi Wang, Lina Wang, Run Wang, Zhibo Wang, and Aoshuang Ye. 2019. Towards a robust deep neural network in texts: A survey. *arXiv preprint arXiv:1902.07285*.
- Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2017a. Ex machina: Personal attacks seen at scale. In *WWW’17*, pages 1391–1399.
- Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2017b. [Wikipedia talk labels: Personal attacks](#).
- Guoyang Zeng, Fanchao Qi, Qianrui Zhou, Tingji Zhang, Bairu Hou, Yuan Zang, Zhiyuan Liu, and Maosong Sun. 2021. [Openattack: An open-source textual adversarial attack toolkit](#). In *ACL’21, Demo*, pages 363–371.

Dataset	#Texts	#Tokens
List of Bad Words ³	1.9K	1.9K
Rumours (Twitter) (Kochkina et al., 2018)	99K	159K
Hate Memes (Twitter) (Gomez et al., 2020)	150K	328K
Personal Atks (Wiki.) (Wulczyn et al., 2017b)	116K	454K
Toxic Comments (Wiki.) (Kaggle, 2019)	2M	1.6M
Malignant Texts (Reddit) (Kaggle, 2021) ⁴	313K	857K
Hateful Comments (Reddit) (Kaggle, 2021) ⁵	1.7M	1M
Sensitive Query (Search Engine, Private)	1.2M	314K
Hateful Comments (Online News, Private)	12.7M	7M
Total texts used to extract ANTHRO	18.3M	-

Table A.1: Real-life datasets that are used to extract adversarial texts in the wild, number of total examples (#Texts) and unique tokens (#Tokens) (case-insensitive)

A Supplementary Materials

Below are list of supplementary materials:

- Table A.1: list of datasets we used to curate the corpus \mathcal{D} , from which human-written perturbations are extracted (Sec. 3.1). All the datasets are publicly available, except from the two private datasets *Sensitive Query* and *Hateful Comments*.
- Table A.2: list of datasets we used to evaluate the attack performance of all attackers (Sec. 4.1) and the prediction performance of BERT and RoBERTa on the respective test sets. All datasets are publicly available.
- Table A.3: Statistical analysis of the human study results (Sec. 4.2).
- Table 4: List of top reasons provided by the professional annotators on why they prefer ANTHRO over *TextBugger* in the Turing test (Sec. 4.2).
- Figure A.1: Word-cloud of extracted human-written perturbations by ANTHRO for some of popular English words.
- Figure A.2, A.3: Interfaces of the human study described in Sec. 4.2.

B Implementation Details

B.1 Attackers

We evaluate all the attack baselines using the open-source *OpenAttack* framework (Zeng et al., 2021). We keep all the default parameters for all the attack methods.

Dataset	#Total	BERT	RoBERTa
CB (Wulczyn et al., 2017a)	449K	0.84	0.84
TC (Kaggle, 2018)	160K	0.85	0.85
HS (Davidson et al.)	25K	0.91	0.97

Table A.2: Evaluation datasets Cyberbullying (CB), Toxic Comments (TC) and Hate Speech (HS) and prediction performance in F1 score on their test sets of BERT and RoBERTa.

Alternative Hypothesis	Mean t-stats	p-value	df
— AMT Workers as Subjects —			
$\mathcal{H}_{\text{Semantics}} : \text{ANTHRO} > \text{TB}$	0.82	5.66	4.1e-7** 48
$\mathcal{H}_{\text{Semantics}} : \text{ANTHRO}_\beta > \text{TB}$	0.64	1.95	2.9e-2* 46
$\mathcal{H}_{\text{Turing}} : \text{ANTHRO} > \text{TB}$	0.71	3.14	1.5e-3** 47
$\mathcal{H}_{\text{Turing}} : \text{ANTHRO}_\beta > \text{TB}$	0.70	3.00	2.2e-3** 46
— Professional Annotators as Subjects —			
$\mathcal{H}_{\text{Semantics}} : \text{ANTHRO} > \text{TB}$	0.75	3.79	2.4e-4** 44
$\mathcal{H}_{\text{Semantics}} : \text{ANTHRO}_\beta > \text{TB}$	0.68	2.49	8.6e-3** 41
$\mathcal{H}_{\text{Turing}} : \text{ANTHRO} > \text{TB}$	0.70	3.06	1.82e-3** 50
$\mathcal{H}_{\text{Turing}} : \text{ANTHRO}_\beta > \text{TB}$	0.73	3.53	4.6e-4** 48
Statistical significant **($p\text{-value} \leq 0.01$) *($p\text{-value} \leq 0.05$)			

Table A.3: It is *statistically significant* ($p\text{-value} \leq 0.01$) that adversarial texts generated by ANTHRO are better than those generated by TextBugger (TB) at both preserving the semantics of the original sentences ($\mathcal{H}_{\text{Semantics}}$) and at being perceived as human-written texts ($\mathcal{H}_{\text{Turing}}$).

B.2 Defenders

For the (1) *Accents normalization*, we adopt the accents removal code from the *Hugging Face* repository⁶. For (2) *Homoglyph normalization*, we adopt a 3rd party python *Homoglyph* library⁷. For (3) *Perturbation normalization*, we use the state-of-the-art misspelling-based perturbation correction *Neuspell* model (Jayanthi et al., 2020)⁸. For *Perspective API*, we directly use the publicly available API provided by Jigsaw and Google⁹.

B.3 Human Study Controls

To ensure a high quality response from MTurks, we require a minimum attentions span of 30 seconds for each question. We recruit MTurks that satisfy the following qualifications, namely (1) recognized as “master” workers by MTurk system, (2) have done at least 5K HITs and (3) have historical HITs approval rate of at least 98%. We pay each worker around \$10 an hour.

⁶ <https://huggingface.co>

⁷ <https://github.com/codebox/homoglyph>

⁸ <https://github.com/neuspell/neuspell>

⁹ <https://www.perspectiveapi.com/>



Read the two pieces of text below and select your response.

Original Text : ==Hello cretino!!!== I'm tired to see too many idiot in Amesh Wiki.

Text A : == Hello cretino !! != I ' m tired to see too many **idior** in Amesh Wiki .

Text B : == Hello cretino !! != I ' m tired to see too many **idipt** in Amesh Wiki .

☐ **Text A** preserves the **meaning** of the Original Text better than **text B**

☐ **Text B** preserves the **meaning** of the Original Text better than **text A**

☐ **Cannot decide (Text A and Text B are similar)**

Submit

Read the two pieces of text below and select your response.

Original Text : He is a CHEATER, and the article should say that.

Text A : He is a CHETAER , and the article should say that .

Text B : He is a CHETER , and the article should say that .

☐ Text A is more likely to be written by a human (and not by a machine) than text B

☐ Text B is more likely to be written by a human (and not by a machine) than text A

☐ Cannot decide (Text A and Text B are similar)

Submit