

AdvExpander: Generating Natural Language Adversarial Examples by Expanding Text

Zhihong Shao¹, Zitao Liu², Jiyong Zhang³, Zhongqin Wu², Minlie Huang^{1*}

¹ Department of Computer Science and Technology, Institute for Artificial Intelligence

¹ State Key Lab of Intelligent Technology and Systems

¹ Beijing National Research Center for Information Science and Technology

¹ Tsinghua University, Beijing 100084, China

² TAL Education Group, Beijing, China ³ Hangzhou Dianzi University

szh19@mails.tsinghua.edu.cn, liuzitao@tal.com

jzhang@hdu.edu.cn, wuzhongqin@tal.com

aihuang@tsinghua.edu.cn

Abstract

Adversarial examples are vital to expose the vulnerability of machine learning models. Despite the success of the most popular substitution-based methods which substitutes some characters or words in the original examples, only substitution is insufficient to uncover all robustness issues of models. In this paper, we present **AdvExpander**, a method that crafts new adversarial examples by expanding text, which is complementary to previous substitution-based methods. We first utilize linguistic rules to determine which constituents to expand and what types of modifiers to expand with. We then expand each constituent by inserting an adversarial modifier searched from a CVAE-based generative model which is pre-trained on a large scale corpus. To search adversarial modifiers, we directly search adversarial latent codes in the latent space without tuning the pre-trained parameters. To ensure that our adversarial examples are label-preserving for text matching, we also constrain the modifications with a heuristic rule. Experiments on three classification tasks verify the effectiveness of AdvExpander and the validity of our adversarial examples. AdvExpander crafts a new type of adversarial examples by text expansion, thereby promising to reveal new robustness issues.

1 Introduction

Adversarial examples are deliberately crafted from original examples to fool machine learning models, which can help (1) reveal systematic biases of data (Zhang et al., 2019b; Gardner et al., 2020), (2) identify pathological inductive biases of models (Feng et al., 2018) (e.g., adopting shallow heuristics (McCoy et al., 2019) which are not robust and unlikely to generalize beyond training data), (3) regularize parameter learning (Minervini and Riedel, 2018),

*Corresponding author: Minlie Huang.

Matched Case Example	
Substitution	<i>Paraphrase</i> (0.754) ✓ → <i>Non-paraphrase</i> (0.731) X
Sentence 1	What are some mind-blowing vehicle accessories that exist that most people don't know about?
Sentence 2	What are some mind-blowing [vehicle → automobile] accessories that exist that most people don't know about?
AdvExpander	<i>Paraphrase</i> (0.754) ✓ → <i>Non-paraphrase</i> (0.707) X
Sentence 1	What are some mind-blowing vehicles tools that exist that most people whom I interviewed don't know about?
Sentence 2	What are some mind-blowing vehicle accessories that exist that most people whom I interviewed don't know about?
Unmatched Case Example	
Substitution	<i>Non-paraphrase</i> (0.991) ✓ → <i>Paraphrase</i> (0.675) X
Sentence 1	What's the best way to send mass emails?
Sentence 2	How can I send mass [emails → emailed] without being aggravating?
AdvExpander	<i>Non-paraphrase</i> (0.991) ✓ → <i>Paraphrase</i> (0.758) X
Sentence 1	What's the best way to send mass emails despite higher fees?
Sentence 2	How can I send mass emails without being aggravating?

Figure 1: Adversarial samples on Quora Question Pairs, crafted against BERT by a substitution-based attack method and the insertion-based AdvExpander. [$A \rightarrow B$] means substituting word A with word B . The underlined expressions are **adversarial modifiers** inserted to expand the **target constituents** in bold.

(4) and evaluate stability (Cheng et al., 2019) or security level of models in practical use.

The most prevalent and effective practice of crafting natural language adversarial examples for classification is to flip characters (Ebrahimi et al., 2018) or substitute words with their typos (Gao et al., 2018; Liang et al., 2018), synonyms (Papernot et al., 2016; Alzantot et al., 2018; Ren et al., 2019; Jin et al., 2019) or other context-compatible words (Zhang et al., 2019a), while reusing the labels of the original examples as long as perturbations are few enough. However, these adversarial attacks limit the search space to the neighborhood of the original text and introduce only small lexical variation, which may not be able to uncover all robustness issues of models.

In this work, we present **AdvExpander**, which crafts new adversarial examples for classification by expanding text under the black-box setting.

Specifically, we **first** use linguistic rules to identify constituents that are safe to expand without leading to an ill-formed text structure. We **then** expand each constituent by inserting an adversarial modifier which is searched from a CVAE-based (Conditional Variational Auto-Encoder (Sohn et al., 2015)) generative model pre-trained on the Billion Word Benchmark (Chelba et al., 2014). We search adversarial modifiers using REINFORCE (Williams, 1992). However, we avoid tuning pre-trained parameters as it can easily sacrifice grammaticality; instead, we additionally introduce a lightweight feed-forward network to search adversaries in the latent space. To make AdvExpander applicable to text matching (e.g., natural language inference and paraphrase identification) besides text classification (e.g., sentiment classification), we design a heuristic rule to ensure modifications are label-preserving: for *matched* cases (e.g., *entailment* pairs in natural language inference and *paraphrase* pairs in paraphrase identification), we only expand shared constituents in both texts with the same modifiers (see the *matched* case in Fig 1).

We characterize AdvExpander in two aspects. **First**, AdvExpander differs from the aforementioned substitution-based attacks considerably. As the semantics of modifiers is far less restricted, and the expressions can be much more diverse than lexical substitutes, AdvExpander has larger search space and can introduce more linguistic variations besides lexical variation, e.g., syntactic variation and semantic variation. Take the *matched* case in Fig 1 for example. For most existing substitution-based attack methods, the candidate substitutes of “*vehicle*” are restricted to its synonyms, e.g., “*automobile*”, “*car*”. By contrast, for AdvExpander, there exist many reasonable modifiers of different types for “*most people*”, e.g., clauses like “*whom I interviewed*”, and prepositional phrases like “*in the neighborhood*”. Therefore, AdvExpander is promising to measure the generalization ability of models. **Second**, as AdvExpander and substitution-based attacks adopt different types of manipulations (ours is based on insertion) and search adversarial examples in different search spaces, they complement each other and can be combined to boost attack performance.

We applied AdvExpander to attack three state-of-the-art models (including RE2 (Yang et al., 2019), BERT (Devlin et al., 2019), and WCNN (Kim, 2014)) and two models with certified robustness

to adversarial word substitutions (Jia et al., 2019) (including bag-of-words and CNN) on SNLI (Bowman et al., 2015), Quora Question Pairs¹, and IMDB² which are commonly used datasets for natural language inference, paraphrase identification, and text classification respectively. We successfully reduce the accuracy of all target models to significantly below-chance level. Furthermore, the validity of our adversarial examples is verified by human evaluation.

Our contributions are summarized as follows: (1) We propose AdvExpander, which generates new adversarial examples by expanding constituents in texts with modifiers. This method is able to introduce rich linguistic variations and differs substantially from existing substitution-based methods; (2) On three classification datasets, AdvExpander substantially degrades the performance of three state-of-the-art models and two models robust to word substitutions, while human annotators remain highly accurate on such adversarial examples, which verifies the validity of our method.

2 Related Work

Adversarial examples are of high value as they can reveal robustness issues of very successful deep classification models. According to how adversarial examples are crafted, recent work can be roughly divided into generation-based ones and edit-based ones.

2.1 Generation-based Adversarial Examples

Some studies utilize rules or neural generation methods to craft adversarial examples. (McCoy et al., 2019) focused on natural language inference and generates an adversarial hypothesis from a premise based on linguistic rules. Though effective, rule-based methods introduce limited variations. (Iyyer et al., 2018) introduced syntactic variation by paraphrasing original text with syntax-controlled network. The generated examples are not optimized to be adversarial. (Kang et al., 2018) utilized Generative Adversarial Nets (Goodfellow et al., 2014) with generator generating adversarial examples and discriminator being the target model. This method is hard to balance grammaticality and adversary. (Zhao et al., 2018) trained an inverter to map a text to a latent representation and

¹<https://data.quora.com/First-Quora-Dataset-Release-QuestionPairs>

²<https://datasets.imdbws.com/>

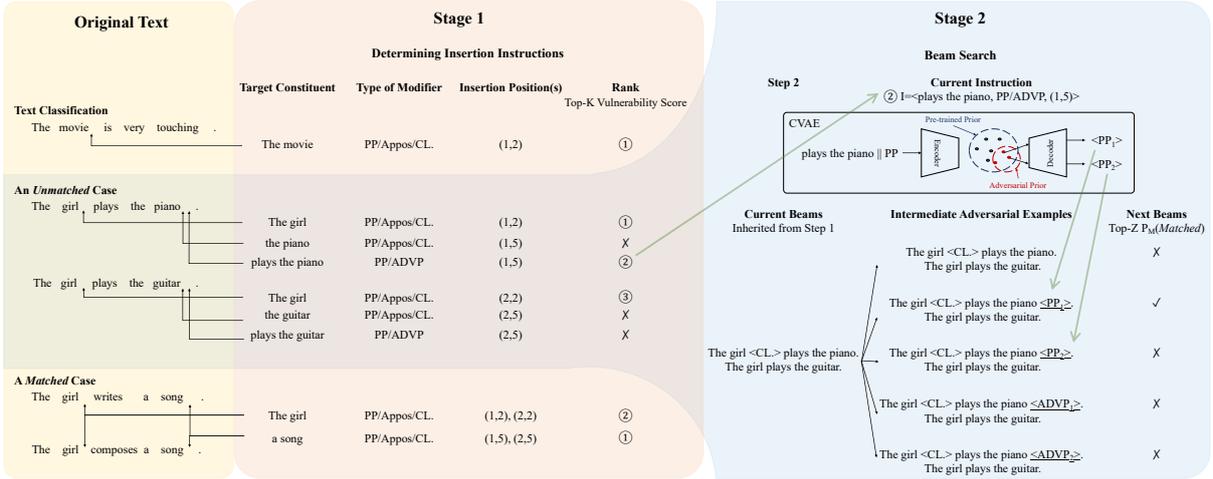


Figure 2: Workflow of AdvExpander for text classification and text matching. “<CL>”, “<PP>”, and “<ADVP>” denote a modifier of type *CL*., *PP*, and *ADVP*, respectively. **Stage 1** mainly determines the insertion instructions, i.e., which constituents to expand, what types of modifiers to expand with, and where to insert the modifiers (we use (i, j) to indicate that the output modifier should be inserted after the j^{th} word of the i^{th} sentence). **Stage 2** is to search adversarial examples via beam search; in the beginning, beams are initialized as the original example. Each beam search step follows one instruction to search adversarial modifiers from a pre-trained CVAE. In this figure, beam size is 1; “<CL>” is the modifier inserted at beam search step 1; underlined modifiers (e.g., “<PP₁>”) are newly inserted modifiers at beam search step 2. If some successful adversarial example(s) is/are found after a beam search step, AdvExpander stops the beam search and returns the successful adversarial example with the lowest perplexity scored by GPT-2.

searched adversaries nearby with heuristic rules. They trained the inverter on the original dataset, which might be insufficient to learn a smooth latent space. Also, their search strategy still has space for improvement. AdvExpander also involves neural generation. By contrast, we choose not to generate a complete text but generate only modifiers which are easier to control and thus less prone to syntactic errors. To learn smooth latent representations of texts, we pre-train a generative model on a large scale corpus. To balance grammaticality, efficiency and effectiveness when finding adversaries, we do not optimize pre-trained parameters but additionally introduce a lightweight feed-forward network to search adversaries in the latent space.

2.2 Edit-based Adversarial Examples

Most studies craft adversarial examples by editing the original text. Substitution is the most popular edit type. Substitution-based attacks are to search an optimal combination of adversarial substitutions under constraints. Under the black-box setting, adversarial attacks often involves scoring the importance of tokens (characters or words), which helps focus attention on important ones to reduce queries and perturbations. A common way of importance scoring is to measure changes of target model out-

put after removing a token (Yang et al., 2018). The optimization process can be conducted by substituting tokens (1) in word order (Papernot et al., 2016), (2) from important ones to less important ones (Ren et al., 2019; Jin et al., 2019; Li et al., 2020a), (3) with beam search (Ebrahimi et al., 2018), (4) or with population-based methods (Alzantot et al., 2018; Zang et al., 2020). Constraints can be grammaticality, semantics-preservation, or context compatibility (Li et al., 2020b).

Compared with substitution which mainly introduces lexical variation, insertion and deletion together are likely to introduce richer variations but are far less popular, as they are more likely to render an adversarial example invalid. (Wallace et al., 2019) inserted universal triggers into text but the inserted strings are meaningless. (Zhang et al., 2019a) supported all three edit types but on token level, thus limited to small perturbations. Most relevantly, (Liang et al., 2018) inserted adversarial phrases which were crafted manually. To the best of our knowledge, AdvExpander is the first efficient method that can automatically insert complex expressions, i.e., modifiers of constituents.

3 Method

4 Task Definition

Suppose a classifier M maps the input text space \mathcal{X} to the label space \mathcal{Y} . Let $X = s_1 s_2 \dots s_N$ be an input text and Y is its label. For text classification, X is a text with N sentences and s_i is the i^{th} sentence. For text matching, $N = 2$ and $\langle s_1, s_2 \rangle$ is the pair to be classified. Our goal is to craft a valid adversarial input X_{adv} by expanding X , so that $M(X_{adv}) \neq Y$. Under the black-box setting, we only have access to the target model’s predictions and the confidence scores.

4.1 Overview

Our method (Fig 2) can be divided into two stages. For convenience, we first introduce the concept “insertion instruction”: an insertion instruction specifies one constituent to expand and the feasible type(s) of modifier to expand with. At **the first stage**, we utilize linguistic rules to analyze the feasible insertion instructions, while ensuring that the insertions will not render the text structure ill-formed. We consider a text structure ill-formed if it is syntactically incorrect or there exists a constituent having multiple modifiers of the same constituency type (e.g., “**The man in white behind the door ...**”). To ensure that insertions are label-preserving for text matching, we further process the instructions so that we only expand shared constituents in both sentences with the same modifiers for *matched* cases (e.g., *entailment* pairs in natural language inference and *paraphrase* pairs in paraphrase identification). For computational efficiency, we only keep the most promising instructions (Eq. 3).

The second stage follows these instructions to find adversarial examples via beam search. At each beam search step, we follow one instruction to search adversarial modifiers in the latent space of a pre-trained generative model. We craft an adversarial example by inserting adversarial modifiers into the original example. In the end, we measure the perplexity of each successful adversarial example with GPT-2³ (Radford et al., 2019), and return the top-ranking one which is expected to be the most syntactically correct.

³https://s3.amazonaws.com/models.huggingface.co/bert/gpt2-torch_model.bin

4.2 Stage 1: Determining Insertion Instructions

AdvExpander expands text by adding modifiers. We consider four types of modifiers, namely adverb phrase (ADVP), prepositional phrase (PP), appositive (Appos), and clause (CL.). For each input sentence s_i , we first obtain its constituency structure⁴, and then utilize handcrafted parsing templates⁵ to determine which constituents to expand and what types of modifiers to expand with. To avoid rendering the text structure ill-formed, we ignore those types of modifiers that the target constituents already have. These analytical results are formatted as a sequence of insertion instructions. Each instruction I is defined as:

$$I = \langle c, t, P \rangle \quad (1)$$

which means a modifier of type t (one of the four types mentioned above) should be inserted into every position within the set P to modify the target constituent c . For example, the *unmatched* case in Fig 2 has three feasible insertion instructions for each sentence. One of the instructions is $I = \langle c = \text{“The girl”}, t = \text{“PP/Appos./CL.”}, P = \{(1,2)\} \rangle$, where $(1,2)$ means that the modifier should be inserted after the 2nd word in the 1st sentence.

To ensure insertions are label-preserving for text matching, we only expand shared constituents in both texts with the exact same modifiers for *matched* cases. Take the *matched* case in Fig 2 for example. We only expand the shared constituents — “The girl” and “a song” — with the exact same modifiers, but ignore the different verb phrases “writes a song” and “composes a song”. Therefore, the instruction associated with “The girl” has insertion position set $P = \{(1,2), (2,2)\}$.

For computational efficiency, we only retain those instructions with top- K vulnerability score:

$$\mathcal{I} = \text{Top-}K_I \text{Score}(I) \quad (2)$$

$$\text{Score}(I) = \max_{X_{adv} \in BS_step(I, \{X\})} 1 - P_M(Y|X_{adv}) \quad (3)$$

where $BS_step(I, \{X\})$ (see the next section) returns S adversarial examples searched in one-step beam search which follows the instruction I and starts from X . $P_M(Y|X_{adv})$ is the probability of Y given the intermediate adversarial example X_{adv} . For an instruction I , vulnerability score measures

⁴<https://s3-us-west-2.amazonaws.com/allennlp/models/elmo-constituency-parser-2018.03.14.tar.gz>

⁵For more details, refer to supplementary material.

the vulnerability of the target constituent by insertion trials. We can calculate vulnerability score for each instruction in parallel.

4.3 Stage 2: Searching Adversarial Examples

The second stage follows the insertion instructions in decreasing order of vulnerability score to search adversarial examples via beam search. During beam search, we maintain a set of intermediate adversarial examples B . At each beam search step, we follow one instruction I , and search adversarial modifiers for each $X_{adv} \in B$ from the latent space of a CVAE-based generative model.

4.3.1 Design of Generative Model

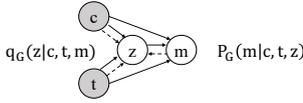


Figure 3: Graphical model of our CVAE-based generative model, where c , t , z , and m denote target constituent, type of output modifier, latent variable, and output modifier, respectively. Dashed arrows are connections for the posterior distribution $q_G(z|c, t, m)$.

Fig 3 shows the graphical model of our generative model. Our generative model takes as input the target constituent c and the expected type of the modifier t , samples the latent variable z , and generates the modifier m . We encode c and decode m with bi- and uni- directional RNNs, respectively.

Our generative model has two critical designs. **First**, we choose CVAE instead of Sequence-to-Sequence (Seq2Seq) model (Bahdanau et al., 2015). This is because Seq2Seq trained with maximum likelihood estimation mainly captures low-level variations of expressions (Serban et al., 2017; Shao et al., 2019), thus failing to provide rich candidates of adversarial modifiers. **Second**, we generate modifiers conditioned on the target constituent instead of the entire text. This makes the distribution of conditions denser, which is beneficial to learn a smoother latent space and to improve the diversity and quality of text. This design also mitigates the gap between pre-trained corpus and attacked corpora, so that we can apply one pre-trained generative model to attack classifiers on different datasets.

4.3.2 Pre-training Generative Model

To learn a smooth latent space, we pre-trained our generative model on the Billion Word Benchmark. For training data, we treated each sentence in this

corpus as having been expanded, and utilize the parsing templates used in stage 1 to extract constituents and their modifiers.

The loss function \mathcal{L}^{pt} is the sum of two terms:

$$\mathcal{L}^{pt} = \mathcal{L}_1^{pt} + \mathcal{L}_2^{pt} \quad (4)$$

The first term \mathcal{L}_1^{pt} is the negative evidence lower bound of $\log P_G(m|c, t)$ which is log likelihood of modifier m given its type t and the modified constituent c :

$$\mathcal{L}_1^{pt} = -E_{q_G(z|c, t, m)}[\log P_G(m|c, t, z)] + D_{KL}(q_G(z|c, t, m)||p_G(z|c, t)) \quad (5)$$

where $p_G(z|c, t)$ and $q_G(z|c, t, m)$ are the prior and posterior distribution of the latent variable z respectively, which are isotropic Gaussian density function. $D_{KL}(\cdot)$ is KL divergence.

The second term \mathcal{L}_2^{pt} is the loss of reconstructing t from the latent variable z , which encodes t into z so that the type of output modifier can be better controlled (Ke et al., 2018).

$$\mathcal{L}_2^{pt} = -E_{q_G(z|c, t, m)}[\log P_G(t|z)] \quad (6)$$

4.4 Beam Search

Before beam search, the set of intermediate adversarial examples is initialized as $B_1 = \{X\}$. Let I_i be the i^{th} promising instruction in \mathcal{I} according to its vulnerability score. The i^{th} beam search step focuses on I_i to expand $X_{adv} \in B_i$. The next beams B_{i+1} is updated as follows:

$$B_{i+1} = \text{Top-}Z_{X'_{adv} \in B_i \cup BS_step(I_i, B_i)}^{1 - P_M(Y|X'_{adv})} \quad (7)$$

where $BS_step(I_i, B_i)$ is a beam search step; following I_i , for each $X_{adv} \in B_i$, it searches S modifiers from the generative model and separately inserts them into X_{adv} , resulting in S new intermediate adversarial examples.

A straightforward way to search adversarial modifiers is to randomly sample S latent codes from the prior distribution and decode a modifier m_{adv} for each code:

$$m_{adv} = \arg \max_m P_G(m|c, t, z), z \sim P_G(z|c, t) \quad (8)$$

As the generative model is capable of producing a rich set of diverse modifiers, this method shows good attack performance. However, this searching process takes no consideration of the target model. We can further optimize the attack performance using REINFORCE(Williams, 1992).

To avoid sacrificing grammaticality for adversary, we choose not to finetune the pre-trained parameters of our generative model but directly search latent codes that will produce adversarial modifiers. As the prior network of our CVAE-based model maps an input to a space of proper but not necessarily adversarial latent codes, we additionally introduce an adversarial prior network, which is a trainable lightweight feed-forward network that narrows the latent space down to adversarial region (see Stage 2 in Fig 2).

The adversarial prior network computes adversarial prior distribution $q_G^{adv}(z|c, t)$ which is isotropic Gaussian density function. The adversarial prior network is initialized with the parameters of the pre-trained prior network, and is finetuned with REINFORCE. The reward is defined as:

$$m_{adv} = \arg \max_m P_G(m|c, t, z), z \sim P_G^{adv}(z|c, t) \quad (9)$$

$$R(z) = -\log(P_M(Y|m_{adv} \rightarrow X_{adv})) + \alpha$$

where m_{adv} is the modifier decoded from the latent code z which is sampled from the adversarial prior distribution, $m_{adv} \rightarrow X_{adv}$ is the intermediate adversarial example crafted by inserting m_{adv} into X_{adv} ($\in B_i$), α is a hyperparameter. To restrict that adversarial latent codes lie in the pre-trained prior distribution and produce modifiers of expected types, we regularize the finetuning procedure with:

$$\Lambda = D_{KL}(q_G^{adv}(z|c, t) || p_G(z|c, t)) - E_{q_G^{adv}(z|c, t)}[\log P_G(t|z)] \quad (10)$$

Therefore, the total loss is:

$$\mathcal{L}^{adv} = -E_{q_G^{adv}(z|c, t)}[R(z)] + \gamma\Lambda \quad (11)$$

The beam search step $BS_step(I_i, B_i)$, for each $X_{adv} \in B_i$, minimizes \mathcal{L}^{adv} for S steps, and returns S new adversarial example (one example per step). Each $X_{adv} \in B_i$ can be processed in parallel.

4.4.1 Finalization of Adversarial Examples

If some intermediate adversarial example $X'_{adv} \in BS_step(I_i, B_i)$ fools the target model, we return the adversarial example with the lowest perplexity:

$$X_{adv}^* = \arg \min_{\substack{X'_{adv} \in BS_step(I_i, B_i) \\ M(X'_{adv}) \neq Y}} perplexity(X'_{adv}) \quad (12)$$

otherwise we continue beam search.

5 Experiments

5.1 Tasks

We evaluated AdvExpander on three datasets: (1) **SNLI**: A large scale dataset for natural language inference which is to judge whether a premise entails, contradicts or is independent of a hypothesis. The train/validation/test split is 550,152/10,000/10,000, respectively. (2) **QQP**: Quora Question Pairs for paraphrase identification which is to identify whether two sentences are paraphrases. The train/validation/test split is 384,348/10,000/10,000, respectively (Wang et al., 2017). (3) **IMDB**: Movie reviews for document-level two-way sentiment classification, with 25,000/25,000 training/test instances respectively.

5.2 Target Models

We attacked RE2 and BERT on both SNLI and QQP, and attacked WCNN and BERT on IMDB. To verify that AdvExpander crafts new adversarial examples, we also attacked two models with certified robustness to adversarial word substitution, i.e., RBOW and RCNN from (Jia et al., 2019).

(1) **RE2**: A simple but effective model which exploits rich alignment features for text matching. (2) **BERT**: Bidirectional Transformer encoder which is pre-trained on large scale corpora. We finetuned BERT_{base-uncased} on the three datasets respectively. (3) **WCNN**: Word-based Convolutional Neural Network. (4) **RBOW**: A robustly trained classifier with bag-of-words encoding. (5) **RCNN**: A robustly trained bag-of-words model with input word vectors transformed by a CNN layer.

5.3 Substitution-based Attack Algorithms

To verify that we can craft new adversarial examples, we compare AdvExpander with three recently proposed black-box substitution-based attack algorithms, i.e., PWWS, BERT-Attack, and TextFooler. The three algorithms mainly differ in their estimation of word importance and the source of substitutes. As the three algorithms are demonstrated to be more effective than or comparable to many other algorithms, they are representative.

PWWS: (Ren et al., 2019) crafts semantic-preserving adversarial examples by replacing words with their synonyms (using WordNet⁶) or replacing named entities with other similar ones. As PWWS has no named entity substitution rules spe-

⁶<https://wordnet.princeton.edu/>

cialized for SNLI or QQP, we applied PWWS on SNLI and QQP without named entity substitutions. **BERT-Attack**: (Li et al., 2020b) crafts adversarial examples by substituting (sub)words with context compatible alternatives sampled from BERT.

TextFooler: (Jin et al., 2019) crafts semantic-preserving adversarial examples and finds synonyms with counter-fitting word embeddings (Mrkšić et al., 2016).

5.4 Implementation Details

Insertion budget K is 3/3/5 for SNLI/QQP/IMDB, respectively. Search steps S is 80 and beam size Z is 5. Thus, AdvExpander queries a target model for no more than 240/240/400 times to craft an adversarial example on SNLI/QQP/IMDB, respectively. **For detailed implementation details, ablation analyses, case study, and error analysis, refer to supplementary material.**

5.5 Automatic Evaluation

Dataset	SNLI			QQP		IMDB		
	RE2	BERT	RBOW	RE2	BERT	WCNN	BERT	RCNN
Ori. Test	86.9	90.7	79.4	88.6	91.3	90.0	92.0	79.3
Adv.	23.8	27.2	21.4	33.0	44.2	10.9	16.4	6.6
Adv. Len.	35.4	36.8	32.7	46.3	50.8	274.3	290.7	268.0
Ori. Len.	23.8	23.8	23.7	24.1	23.8	245.0	257.1	238.0
Ori. Test (long)	86.6	89.9	77.3	95.0	96.4	89.3	87.9	78.8

Table 1: Automatic evaluation performance, including model accuracy on the original test examples (“*Ori. Test*”) and model accuracy on the corresponding adversarial examples crafted by AdvExpander (“*Adv.*”). “*Adv. Len.*” and “*Ori. Len.*” denote the average length of **successful** adversarial examples and the **corresponding** original examples before expansion, respectively. “*Ori. Test (long)*” denotes model accuracy on the original test examples that are longer than the average length of successful adversarial examples.

We evaluated AdvExpander on the entire test sets for SNLI and QQP but on 1,000 random test samples for IMDB, as texts in IMDB are hundreds of words long and even the baseline PWWS is slow. We measured model accuracy on the original test examples and the corresponding adversarial examples respectively (Table 1).

AdvExpander is effective in crafting adversarial examples; it degrades the accuracy of all target models substantially. For example, the accuracy of BERT drops from above 90% to below 28% on both SNLI and IMDB.

As AdvExpander crafts adversarial examples by expanding text, we further investigate the influence of text length on model accuracy. As shown in

Dataset	SNLI			QQP		IMDB		
	RE2	BERT	RBOW	RE2	BERT	WCNN	BERT	RCNN
PS	28.7	36.3	41.2	49.2	53.5	7.9	29.1	14.8
BA	10.1	12.8	17.6	30.3	36.6	0.8	12.1	9.0
TF	2.9	4.4	8.9	32.1	38.1	0.4	14.3	10.0
PS + BA	5.5*	8.1*	11.4*	29.6*	35.6*	0.6*	11.0*	3.0*
PS + TF	1.6*	3.0*	6.4*	31.4*	37.0*	0.1	13.1*	4.1*
BA + TF	1.5*	2.3*	4.5*	26.6*	32.8*	0.0	10.5*	2.7
Ours + PS	3.6*	5.5*	10.1*	21.0*	29.2*	2.9*	12.9*	2.6*
Ours + BA	1.6*	1.9*	5.6*	15.5	22.7	0.4*	7.2	1.9
Ours + TF	0.3	0.4	1.3	16.4*	23.8*	0.2	8.6*	1.8

Table 2: Comparison between AdvExpander and substitution-based attack methods. “*PS*”/“*BA*”/“*TF*” denotes PWWS/BERT-Attack/TextFooler, respectively. All numbers are model accuracy on the adversarial examples crafted by different attack methods. We also report model accuracy under the attacks of any two combined methods (e.g., “*PS + BA*”): an attack is successful if at least one algorithm fools the target model. Accuracy that is significantly higher than the lowest accuracy (in bold) is marked with * for p-value < 0.05 according to bootstrap resampling (Koehn, 2004).

Table 1, though AdvExpander makes input texts longer, the target models remain high accuracy on the original test examples that are longer than the average length of adversarial examples, indicating that text length is not the factor why AdvExpander is successful to fool target models.

Comparison with Substitution-based Attacks

We further investigate the relationship between AdvExpander and previous substitution-based attack algorithms (Table 2).

AdvExpander degrades model accuracy more remarkably than PWWS in the most cases, but less remarkably than BERT-Attack and TextFooler. Note that we choose not to modify an example if any insertion will render the text structure ill-formed. When ignoring those examples we choose not to modify, the accuracy of BERT under our attacks is 6.2% on SNLI and 37.1% on QQP, respectively, which is competitive with the performance of BERT-Attack and TextFooler.

To verify that AdvExpander crafts new adversarial examples compared with substitution-based methods, we attacked RBOW and RCNN which have certified robustness to word substitutions. Due to robust training, the two models are even harder to fool than some structurally more advanced models for substitution-based attacks. However, as RBOW and RCNN have rather simple architecture and are only trained to be robust to word substitutions, they are unsurprisingly easier to fool than the other models for AdvExpander. Take IMDB for example. RCNN is significantly more accurate (10.0%) than

WCNN (0.4%) under TextFooler’s attacks, but is much less accurate (6.6%) than WCNN (10.9%) under our attacks. Therefore, certified robustness to word substitutions may not indicate robustness to insertion-based adversarial examples.

We also combined AdvExpander with a substitution-based method to attack the target models (Table 2). Specifically, an attack is considered successful if at least one of the two methods fools the target model. The combinatorial attacks consistently boost attack performance. In the most cases, the highest performance boost is brought by combining AdvExpander with a substitution-based method but not by combining two substitution-based methods. In other words, AdvExpander can craft adversarial examples in a way substitution-based methods is incapable of. Thus, AdvExpander is complementary to substitution-based methods and is promising to reveal new robustness issues.

5.6 Human Evaluation

Dataset	QQP		IMDB	
Metrics	Accuracy	Grammaticality	Accuracy	Naturalness
Ori. Sampled	85.0	2.65	88.0	2.69
TF	71.5	2.45	84.0	2.57
Ours	80.0**	2.39	84.5	2.65*

Table 3: Human evaluation of adversarial examples against BERT, in terms of human accuracy and grammaticality/naturalness. “TF” denotes TextFooler. “Ori. Sampled” shows evaluation on the corresponding original test samples. Bootstrap resampling (Koehn, 2004) is used as significance test between the two methods. ** marks significantly better performance for p-value \leq 0.01, and * for p-value $<$ 0.05.

To verify the validity of our adversarial examples, we conducted human evaluation (Table 3). We randomly sampled 200 adversarial examples against BERT on QQP and IMDB, respectively. These samples are mixed with the corresponding original test samples and the corresponding adversarial examples crafted by TextFooler; each example is presented to three workers on Amazon Mechanical Turk to annotate its label and whether it is grammatical/natural⁷ (3-point Likert Scale). For each example, we aggregated human-predicted labels with majority vote, and computed human accuracy as the consistency between the aggregated labels and the gold labels. We also computed the average grammaticality/naturalness score.

⁷As IMDB reviews are informal and contain grammatical errors, we measure naturalness on IMDB.

Human accuracy on the original examples and our adversarial examples is close. By contrast, on TextFooler’s adversarial examples, human accuracy drops to 71.5% on QQP, mostly due to imperfection of synonym candidates (e.g., substituting “mechanical” in “mechanical engineer” with “mechanised”). Therefore, our adversarial examples are label-preserving at an acceptable level. Moreover, the grammaticality/naturalness score of our adversarial examples is close to that of the original samples, indicating that our adversarial examples are of good quality. Overall, these results demonstrate the validity of our adversarial examples.

5.7 Adversarial Training

We separately retrained RE2 on SNLI augmented with 80K adversarial examples crafted on the training set by AdvExpander and TextFooler, and tested their robustness on the original test set (Table 4).

For both TextFooler and AdvExpander, adversarial training helps improve a model’s robustness to the attack method it is trained with, and slightly improves model accuracy on the original test set. Notably, as the original training set is large, training models with more adversarial examples can further improve models’ robustness.

We also observed that adversarially training RE2 with TextFooler can hardly improve accuracy under AdvExpander’s attacks (23.8% \rightarrow 24.1%), and vice versa (2.9% \rightarrow 3.1%). After augmenting the training set with both AdvExpander’s and TextFooler’s adversarial examples (80K for each), we improved model accuracy under AdvExpander’s attacks (23.8% \rightarrow 30.8%) and TextFooler’s attacks (2.9% \rightarrow 7.4%). This indicates that AdvExpander can generate new adversarial examples, and can reveal robustness issues that TextFooler fails to reveal. They complement each other.

Training Set	Ori. Train	+ TF	+ Ours	+ Ours & TF
Ori. Test	86.9	87.3 (+0.4)	87.0 (+0.1)	87.2 (+0.3)
TF	2.9	8.0 (+5.1)	3.1 (+0.2)	7.4 (+4.5)
Ours	23.8	24.1 (+0.3)	30.0 (+6.2)	30.8 (+7.0)

Table 4: Model accuracy on the original test examples (“Ori. Test”) and adversarial examples (“TF” and “Ours”) after adversarially training RE2 on SNLI with TextFooler (“+TF”), with AdvExpander (“+Ours”), or with both AdvExpander and TextFooler (“+Ours & TF”). “TF” stands for TextFooler. Numbers in parentheses are improvements over the model trained on the original training set (“Ori. Train”).

6 Conclusion

In this paper, we present AdvExpander which generates new natural language adversarial examples by expanding text. Extensive experiments demonstrate the effectiveness of our algorithm and the validity of our adversarial examples. Our adversarial examples are substantially different from previous substitution-based adversarial examples, thus promising to reveal new robustness issues.

References

- Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani B. Srivastava, and Kai-Wei Chang. 2018. [Generating natural language adversarial examples](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2890–2896. Association for Computational Linguistics.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. [Neural machine translation by jointly learning to align and translate](#). In *International Conference on Learning Representations*.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 632–642. The Association for Computational Linguistics.
- Jill Burstein, Christy Doran, and Tamar Solorio, editors. 2019. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics.
- Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2014. [One billion word benchmark for measuring progress in statistical language modeling](#). In *INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association, Singapore, September 14-18, 2014*, pages 2635–2639. ISCA.
- Yong Cheng, Lu Jiang, and Wolfgang Macherey. 2019. [Robust neural machine translation with doubly adversarial inputs](#). In (Korhonen et al., 2019), pages 4324–4333.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In (Burstein et al., 2019), pages 4171–4186.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. [Hotflip: White-box adversarial examples for text classification](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers*, pages 31–36. Association for Computational Linguistics.
- Shi Feng, Eric Wallace, Alvin Grissom II, Mohit Iyyer, Pedro Rodriguez, and Jordan Boyd-Graber. 2018. [Pathologies of neural models make interpretations difficult](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3719–3728, Brussels, Belgium. Association for Computational Linguistics.
- Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun Qi. 2018. [Black-box generation of adversarial text sequences to evade deep learning classifiers](#). In *2018 IEEE Security and Privacy Workshops, SP Workshops 2018, San Francisco, CA, USA, May 24, 2018*, pages 50–56. IEEE Computer Society.
- Matt Gardner, Yoav Artzi, Victoria Basmova, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hanna Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. [Evaluating NLP models via contrast sets](#). *CoRR*, abs/2004.02709.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. 2014. [Generative adversarial nets](#). In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680.
- Iryna Gurevych and Yusuke Miyao, editors. 2018. *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*. Association for Computational Linguistics.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. [Adversarial example generation with syntactically controlled paraphrase networks](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 1875–1885. Association for Computational Linguistics.
- Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. 2019. [Certified robustness to adversarial word substitutions](#). In *Proceedings of the*

- 2019 *Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 4127–4140. Association for Computational Linguistics.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2019. [Is BERT really robust? natural language attack on text classification and entailment](#). *CoRR*, abs/1907.11932.
- Dongyeop Kang, Tushar Khot, Ashish Sabharwal, and Eduard H. Hovy. 2018. [Adventure: Adversarial training for textual entailment with knowledge-guided examples](#). In (Gurevych and Miyao, 2018), pages 2418–2428.
- Pei Ke, Jian Guan, Minlie Huang, and Xiaoyan Zhu. 2018. [Generating informative responses with controlled sentence function](#). In (Gurevych and Miyao, 2018), pages 1499–1508.
- Yoon Kim. 2014. [Convolutional neural networks for sentence classification](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1746–1751. ACL.
- Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *EMNLP*, pages 388–395. ACL.
- Anna Korhonen, David R. Traum, and Lluís Màrquez, editors. 2019. *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*. Association for Computational Linguistics.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020a. [BERT-ATTACK: adversarial attack against BERT using BERT](#). *CoRR*, abs/2004.09984.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020b. [BERT-ATTACK: Adversarial attack against BERT using BERT](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202. Association for Computational Linguistics.
- Bin Liang, Hongcheng Li, Miaoqiang Su, Pan Bian, Xirong Li, and Wenchang Shi. 2018. [Deep text classification can be fooled](#). In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pages 4208–4215. ijcai.org.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. [Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference](#). In (Korhonen et al., 2019), pages 3428–3448.
- Pasquale Minervini and Sebastian Riedel. 2018. [Adversarially regularising neural NLI models to integrate logical background knowledge](#). In *Proceedings of the 22nd Conference on Computational Natural Language Learning, CoNLL 2018, Brussels, Belgium, October 31 - November 1, 2018*, pages 65–74. Association for Computational Linguistics.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Blaise Thomson, Milica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, David Vandyke, Tsung-Hsien Wen, and Steve Young. 2016. [Counter-fitting word vectors to linguistic constraints](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 142–148, San Diego, California. Association for Computational Linguistics.
- Nicolas Papernot, Patrick D. McDaniel, Ananthram Swami, and Richard E. Harang. 2016. [Crafting adversarial input sequences for recurrent neural networks](#). In *2016 IEEE Military Communications Conference, MILCOM 2016, Baltimore, MD, USA, November 1-3, 2016*, pages 49–54. IEEE.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8).
- Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. [Generating natural language adversarial examples through probability weighted word saliency](#). In (Korhonen et al., 2019), pages 1085–1097.
- Iulian Vlad Serban, Alessandro Sordani, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron C. Courville, and Yoshua Bengio. 2017. [A hierarchical latent variable encoder-decoder model for generating dialogues](#). In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 3295–3301. AAAI Press.
- Zhihong Shao, Minlie Huang, Jiangtao Wen, Wenfei Xu, and Xiaoyan Zhu. 2019. [Long and diverse text generation with planning-based hierarchical variational model](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 3255–3266. Association for Computational Linguistics.
- Kihyuk Sohn, Honglak Lee, and Xinchen Yan. 2015. Learning structured output representation using deep conditional generative models. In *NIPS*, pages 3483–3491.
- Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. 2019. [Universal adversarial triggers for attacking and analyzing nlp](#). *Proceedings of*

the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP).

Zhiguo Wang, Wael Hamza, and Radu Florian. 2017. [Bilateral multi-perspective matching for natural language sentences](#). In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pages 4144–4150. ijcai.org.

Ronald J. Williams. 1992. [Simple statistical gradient-following algorithms for connectionist reinforcement learning](#). *Machine Learning*, 8:229–256.

Puyudi Yang, Jianbo Chen, Cho-Jui Hsieh, Jane-Ling Wang, and Michael I. Jordan. 2018. [Greedy attack and gumbel attack: Generating adversarial examples for discrete data](#). *CoRR*, abs/1805.12316.

Runqi Yang, Jianhai Zhang, Xing Gao, Feng Ji, and Haiqing Chen. 2019. [Simple and effective text matching with richer alignment features](#). In (Korhonen et al., 2019), pages 4699–4709.

Yuan Zang, Fanchao Qi, Chenghao Yang, Zhiyuan Liu, Meng Zhang, Qun Liu, and Maosong Sun. 2020. [Word-level textual adversarial attacking as combinatorial optimization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 6066–6080. Association for Computational Linguistics.

Huangzhao Zhang, Hao Zhou, Ning Miao, and Lei Li. 2019a. [Generating fluent adversarial examples for natural languages](#). In (Korhonen et al., 2019), pages 5564–5569.

Yuan Zhang, Jason Baldridge, and Luheng He. 2019b. [PAWS: paraphrase adversaries from word scrambling](#). In (Burststein et al., 2019), pages 1298–1308.

Zhengli Zhao, Dheeru Dua, and Sameer Singh. 2018. [Generating natural adversarial examples](#). In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.