

# Lost In Translation: Generating Adversarial Examples Robust to Round-Trip Translation

Neel Bhandari

RV College of Engineering  
neelbhandari64@gmail.com

Pin-Yu Chen

IBM Research  
pin-yu.chen@ibm.com

## Abstract

Language Models today provide a high accuracy across a large number of downstream tasks. However, they remain susceptible to adversarial attacks, particularly against those where the adversarial examples maintain considerable similarity to the original text. Given the multilingual nature of text, the effectiveness of adversarial examples across translations and how machine translations can improve the robustness of adversarial examples remain largely unexplored. In this paper, we present a comprehensive study on the robustness of current text adversarial attacks to round-trip translation. We demonstrate that 6 state-of-the-art text-based adversarial attacks do not maintain their efficacy after round-trip translation. Furthermore, we introduce an intervention-based solution<sup>1</sup> to this problem, by integrating Machine Translation into the process of adversarial example generation and demonstrating increased robustness to round-trip translation. Our results indicate that finding adversarial examples robust to translation can help identify the insufficiency of language models that is common across languages, and motivate further research into multilingual adversarial attacks.

## 1 Introduction

Language models, despite their remarkable success across tasks, have shown to be vulnerable to adversarial examples, which are inputs designed to be similar to the model’s native data inputs, but crafted with small modifications to fool the model during inference. These examples can be classified correctly by a human observer, but often mislead a target model, providing an insight into their robustness to adversarial inputs (Chen and Liu, 2023;

Chen and Hsieh, 2023). They are essential in understanding key vulnerabilities in models across a variety of applications (Chen and Das, 2023).

ML models are being increasingly deployed commercially for translation. A special form of translation is round trip translation, which focuses on translating a given text from one language to the second and back to the first. Round trip translation has been increasingly used in several research areas, including correcting grammatical errors (Lichtarge et al., 2019; Madnani et al., 2012), evaluating machine translation models (Crone et al., 2021; Cao et al., 2020; Moon et al., 2020), paraphrasing (Guo et al., 2021) and rewriting questions (Chu et al., 2020). It is also used extensively as part of the quality assurance process in critical domains such as medical, legal and market search domains. The use of ML models in these critical domains means that they have to be tested by robust adversarial attacks to make for safe and reliable commercial deployment. Given the importance of round trip translation, we are motivated to study its effects on current adversarial attacks.

We summarise our contributions as follows:

- We demonstrate that round trip translation can be used as a cheap and effective defence against *current* textual adversarial attacks. We show that 6 state-of-the-art adversarial text attacks suffer an average performance loss of 66%, rendering most examples generated non-adversarial.
- However, we find that round-trip translation defensive capabilities can be bypassed by our proposed *attack-agnostic algorithm* that provides machine translation intervention to increase robustness against round-trip translation. We find it provides minimal difference in quantification metrics to the original, which shows our method finds a new set of robust and high-quality text adversarial examples

<sup>1</sup>Code for the paper: [https://github.com/neelbhandari6/NMT\\_Text\\_Attack](https://github.com/neelbhandari6/NMT_Text_Attack).  
Emails: Neel Bhandari: neelbhandari64@gmail.com  
Pin-Yu Chen: pin-yu.chen@ibm.com

against neural machine translation (NMT).

## 2 Related Works

(Papernot et al., 2017) proposed a white box adversarial attack that repeatedly modified the input text till the generated text fooled the classifier. This method, although effective in principle, did not maintain semantic meaning of the sentence. (Ebrahimi et al., 2018) and (Samanta and Mehta, 2017) proposed gradient-based solutions involving token based changes and searching for important words. These methods, however, did not prove to be scalable and lacked robust performance. It was followed by methods such as character replacement (Ribeiro et al., 2018), phrase replacement and word scrambling. These techniques, however, fail to maintain semantic consistency with the original input. (Jia et al., 2019) introduced adding distracting sentences to the reading comprehension task. (Jin et al., 2020) propose TextFooler which generates adversaries using token-level similarity and is bound by axiomatic constraints. (Lei et al., 2019) propose paraphrasing attacks using discrete optimization. (Garg and Ramakrishnan, 2020) introduce BAE, which uses masked-language modelling to generate natural adversarial examples for the text. Recent works in adversarial attacks on NMT include (Cheng et al., 2019) using gradient based adversarial inputs to improve robustness of NMT models, and (Zhang et al., 2021) proposed a novel black-box attack algorithm for NMT systems. However, none of these works target round-trip translation, and do not demonstrate attack agnostic capabilities.

## 3 NMT-Text-Attack

In order to generate adversarial examples robust to round-trip translation, we propose an intervention-based attack-agnostic method that only requires access to a neural machine translation(NMT) model, shown in Algorithm 1. We employ a generic template used by standard state-of-the-art adversarial attack examples in order to showcase the attack-agnostic capabilities. From (Li et al., 2019; Jin et al., 2020; Ren et al., 2019; Garg and Ramakrishnan, 2020; Gao et al., 2018) it can be seen that the attacks follow a two section split. The first section is word importance ranking, and the second section deals with word replacement and constraint evalua-

---

### Algorithm 1: NMT-Text-Attack

---

**Input** : Sentence  $S = [w_1, w_2, \dots, w_n]$ ,  
Ground truth label  $Y$ , Victim  
Model  $V$ , Machine Translation  
model  $M$ , User-Specific  
Constraints  $C$ , Attack  $A$

**Output** : Adversarial Example  $X_{adv}$

---

```

1 Phase I - Word Importance Ranking
2 Call attack  $A$ 
3 Initialize edge weights
4 for each word  $w_i$  in  $S$  do
5   | Compute Importance score  $I_i$  from  $A$ 
6 Sort words in descending order into list  $W$ 
7 Phase 2 - Word Replacement
8 # Word Replacement Strategy
9 for each word  $w_i$  in  $W$  do
10  | Predict Top-K replacements for  $w_i$ 
    | using  $A$  and store in  $R = [r_1, r_2, \dots, r_k]$ 
11  for each word  $w_i$  in  $W$  do
12    | Replace  $w_i$  with  $r_j$  in  $S$  to make
    |  $X_{adv}$ 
13    | Round-Trip-Translate  $X_{adv}$  with  $k$ 
    | language(s) using  $M$  to make
    |  $T = [t_1, t_2, \dots, t_p]$  where  $t_i$  is  $X_{adv}$ 
    | translated through language  $i$ 
14    | Evaluate classification scores for
    |  $T = [t_1, t_2, \dots, t_p]$  using  $V$ ,
    | removing examples that do not
    | maintain adversarial sentiment
15    | for each  $c_i \in C$  do
16      | Apply constraint  $c_i$  to each
    |  $t_i \in T$ 
17    | Select best  $t_i \in T$  w.r.t constraints
    |  $C$  and store as  $X_{adv}$ 
18 return  $X_{adv}$ 

```

---

tion, where NMT-Text-Attack is introduced along with the original algorithm’s constraints.

**I. Word Importance Selection.** This section initially involves pre-processing the input sentence with techniques such as removing stop words etc. This is followed by analysing the most important keywords in the target sentence using several techniques, ranging from the input deletion method, to probability weighted word saliency. These methods are specific to the adversarial attack chosen to be integrated with NMT-Text-Attack. For example, TextFooler uses the input deletion method. Once the most important words are learnt, attack algorithms look for replacements through synonym

search or by replacing individual characters of the original input word to make an adversarial candidate.

**II. Constraint Evaluation.** We introduce the machine translation task in this section. First, we predict the Top-K replacements for each word  $w_i$  in word importance ranking list  $W$  and substitute them in the sentence  $S$  iteratively (Step 12). We then implement round-trip translation on these sentences for  $k$  languages, where  $k$  is specified by the user (Step 13). On collecting the candidate sentences, we evaluate them on the sentiment classification model  $V$  and remove all examples that do not maintain the adversarial sentiment post round-trip translation (Step 14). Finally, we apply the algorithm-specific constraints  $C$  on the collected final sentences  $T$ , and select the best candidate based on their similarity score with respect to the original sentence.

This is followed by applying algorithm-specific constraints  $C$  such as semantic similarity to original input on replacement, POS tag preservation etc.

## 4 Evaluation

For performance evaluation, we consider using a range of algorithms from the TextAttack library (Morris et al., 2020).

### 4.1 Dataset and Victim Model

We use the Rotten Tomatoes Movie Reviews and Yelp Polarity datasets to perform sentiment analysis. We sample 1000 random examples from the test set of each of these mentioned datasets and run our experiments on them. For our Victim Model, we use the Bidirectional Encoder Representations from Transformers (BERT) model (Devlin et al., 2019).

### 4.2 Current Attacks are not Robust to Round Trip Translation

We run 6 adversarial attacks on the Movie Reviews Dataset and analyse their robustness to round-trip translation, as shown in Figure 1. We analyse them against 3 languages – Spanish, German and French through the EasyNMT library (see Appendix for more details). On round-trip translating the adversarial examples, we test the resultant examples against the classification model.

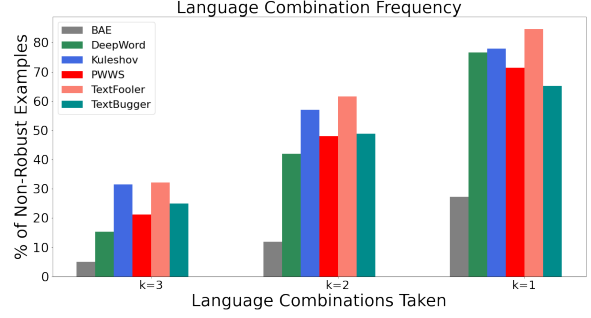


Figure 1: Percentage of non-robust examples flagged by at least  $k$  language combination

On the y-axis, we provide the percentage of non-robust examples to at least  $k$  out of  $m = 3$  languages. Formally, if  $k$  is the number of languages used in tandem,  $N$  is the number of examples in total,  $y_a$  is the original prediction before round trip translation and  $\hat{y}_a$  is the prediction after round-trip translation by translation model  $M$  and victim model  $V$ , then the y-axis is defined as  $Y = \frac{1}{N} \sum_{a=1}^N \mathbb{1}\{\text{at least } k \text{ languages have } y_a \neq \hat{y}_a\}$ , where  $\mathbb{1}\{E\}$  is an indicator function such that it is one when the event  $E$  is true and zero otherwise.

We see that on average, over 66% of the examples generated originally by the attack are rendered non-adversarial on round-trip translation with at least one language ( $k = 1$ ). BAE remains the most robust to translations, while TextFooler remains the least robust. On increasing the number of language combinations taken ( $k > 1$ ), we see that there is a decrease in effectiveness of round trip translation as a defense against the adversarial examples, however there is still significant loss in attack success rate. This is because when you add more languages as a constraint, there is an increased chance that at least one of the constrained languages is robust to round-trip translation for any example. This provides considerable evidence that round trip translation can be used as a cheap and effective defense, and motivates the question of whether there exists text adversarial examples robust to round-trip translation. In the following sections, We evaluate the robustness of our proposed NMT-Text-Attack as shown in Algorithm 1.

### 4.3 NMT-Text-Attack Results

We analyse the results of incorporating NMT-Text-Attack into existing attacks across the mentioned datasets. We evaluate the attack on its success rate with respect to the attacks’ native success rate

Table 1: Success Rate (%) of NMT-Text-Attack Relative to when Original Attack Success Rate is 100% (Replacement generation limit=40)

Dataset	TextFooler+NMT	TextBugger+NMT	PWWS+NMT
MR	70.7	74.7	69.4
Yelp	60.0	71.4	68.8

Table 2: Sentence similarity analysis on Yelp and Movie Reviews (MR) Datasets

Dataset	Attack	USE	Jaccard	BERT
Yelp	TextBugger	0.93	0.79	0.95
	TextFooler	0.93	0.81	0.97
	PWWS	0.93	0.85	0.97
	TextBugger + NMT	0.94	0.848	0.9715
	TextFooler + NMT	0.82	0.724	0.956
	PWWS + NMT	0.83	0.645	0.9265
MR	TextBugger	0.93	0.79	0.95
	TextFooler	0.813	0.715	0.953
	PWWS	0.85	0.77	0.96
	TextBugger + NMT	0.91	0.68	0.92
	TextFooler + NMT	0.82	0.724	0.956
	PWWS + NMT	0.83	0.645	0.9256

without NMT-Text-Attack. Note that, through our novel intervention-based algorithm, we are able to guarantee 100% robustness to back-translation on the user’s selected language(s). This is because our algorithm (line 14) introduces a strict constraint to only allow examples that are robust to back-translation to be selected as candidates for the attack, which leads to significant increase over the original algorithm’s robustness to round-trip translation. This guarantee is important as it helps achieve high-quality robustness in multilingual settings, which no existing adversarial attack can provide. Table 1 shows that to meet this criteria, NMT-Text-Attack is successful on average 30% less examples than it’s original counterpart.

While this loss may seem significant, we believe this is justified for two reasons. First, this loss comes with a 100% success in robustness to round-trip translation coupled with attack success. This is critical in commercial settings where deployed models need to have confident outputs in the face of several language translations. Secondly, in Figure 2, we see that there is considerable scope to increase the number of robust examples available simply by increasing the replacement limit. We set our replacement limit at 40 for our experiments, and Figure 2 demonstrates that scaling the number of replacements significantly increases number of available robust examples.

We also provide a quantitative analysis of our model by analysing the adversarial examples gener-

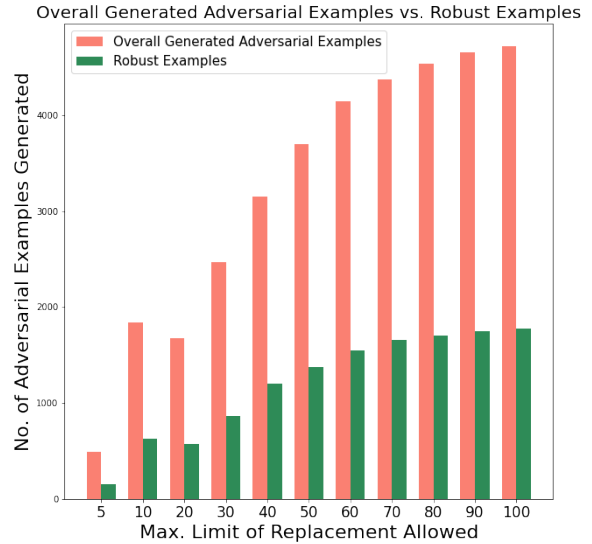


Figure 2: Replacement vs. Robust Examples

ated against the original attack in Table 2. Universal Sentence Encoder (Cer et al., 2018) with cosine similarity, along with Jaccard Similarity are used as similarity metrics, while BERT Score (Zhang et al., 2020) is used to analyse meaning preservation. We notice that there is little variation in the effectiveness of the algorithms when it comes to meaning preservation and similarity, which shows that our proposed intervention, while increasing robustness significantly, maintains the quality of the original attack. Examples of adversarial examples on sentences have been mentioned in the Appendix.

#### 4.4 Ablation Study

In this section, we provide an ablation study to substantiate the performance of our algorithm. In this study, we provide TextFooler with NMT-Text-Attack with 2 ‘seen’ languages and test its performance with an ‘unseen’ language. A ‘seen’ language is defined as one which model is provided with as constraints for adversarial examples to satisfy, as shown in Algorithm 1. An ‘unseen’ language, consequently, is one which the model has not added as a constraint, hence does not guarantee 100% robustness against. The three languages we use are French, German, and Spanish. We alternate between using two of the languages as ‘seen’, and one as ‘unseen’. We compare this with the performance of TextFooler without NMT-Text-Attack on the unseen languages in Table 3. We observe that TextFooler with NMT-Text-Attack outperforms TextFooler without NMT-Text-Attack on



Table 3: Performance of NMT-Text-Attack on unseen language

Seen Languages	Unseen Language	TextFooler +NMT	TextFooler w/o NMT
French and German	Spanish	72.9%	50.61
French and Spanish	German	74.08%	51.97
German and Spanish	French	67%	50.8

Table 4: BLEU and % words perturbed results of NMT-Text-Attack on Yelp and Movie Reviews(MR) Datasets

Dataset	Algorithm	BLEU Score	% Words Perturbed
MR	TextFooler	0.37	16.07
	TextFooler + NMT	0.48	19.33
	TextBugger	0.47	5.17
	TextBugger + NMT	0.62	11.75
	PWWS	0.43	11.57
Yelp	PWWS + NMT	0.57	15.19
	TextFooler	0.50	41.43
	TextFooler + NMT	0.68	56.27
	TextBugger	0.50	34.56
	TextBugger + NMT	0.53	34.60
	PWWS	0.53	35.90
	PWWS + NMT	0.73	50.91

average by 20%. This shows that the integration of our attack-agnostic algorithm provides significant performance increase even in situations where the attack is facing unseen languages.

To provide further substantiate the performance of NMT-Text-Attack, we provide a detailed set of results in Table 4. Here, we see that algorithms with NMT-Text-Attack consistently provide higher BLEU scores than their original versions by a significant margin. We see that the percentage of words perturbed remains lower in the original algorithm. However, given the combination of higher performance across the mentioned metrics and generalisation to unseen languages, we believe that this result justifies itself.

## 5 Conclusion

In this paper, we demonstrate the ineffectiveness of current text adversarial attack algorithms to round-trip translation, and provide an intervention-based method to improve robustness to round-trip translation in these algorithms. We show that this intervention (NMT-Text-Attack) has minimal effect on the actual semantic metrics but can significantly improve the attack success rate against back-translation, suggesting that there exist a new set of robust text adversarial examples. The attack-agnostic nature of the algorithm along with its high-quality performance makes it an effective error diagnosing tool with any existing text attack for inspecting model robustness.

## 6 Appendix

### 6.1 Ethical Concerns

Our paper discusses the potential weakness of NLP models to round-trip translation, and describes an algorithm that can make the weakness more robust. However, we believe that we give new insights in studying text adversarial examples and will spur more robust machine learning models in the future. We are also the first individuals to introduce the vulnerability to round-trip translation, which provides opportunity to develop robust models in a novel setting.

### 6.2 Computational Resources

For the implementation of our algorithm and experiments, we use Google Colab as our base GPU provider. The GPU typically provided is Tesla P100. We use 190 GPU hours to run all our experiments. We use a pre-trained BERT model with 12-head attention and 110 million parameters, which is typical of BERT models.

### 6.3 Machine Translation Setup

We use the Opus-MT set of models through the EasyNMT library (Tang et al., 2020). Opus-MT consists of 1200 models trained on several languages for open translation. The architecture for the Opus-MT models is based on a standard transformer setup with 6 self-attentive layers in both, the encoder and decoder network with 8 attention heads in each layer. This architecture is used to back-translate the target reviews from English to French, German and Spanish, and back to English.

### 6.4 Adversarial Attack Settings

Algorithm 1 details a general template of several state of the art adversarial attacks we have used in the paper. In this section we detail the exact settings used for each adversarial attack when integrated with NMT-Text-Attack. These are standard approaches used directly from the TextAttack Library with no changes in standard settings.

#### 6.4.1 Textfooler

- Word Importance Selection
  - Max allowable replacement candidate generation for synonyms: 40.

- Transformation Embedding Mechanism: Counterfitted Glove Embeddings ([Mrkšić et al., 2016a](#))

- Word Replacement:

- Pre-transformation constraints:
  - \* RepeatModification: A constraint disallowing the modification of words which have already been modified
  - \* StopwordModification: A constraint disallowing the modification of stopwords
- Constraints:
  - \* Minimum cosine distance between word embeddings = 0.5
  - \* Part of Speech : Only replace words with the same part of speech (or nouns with verbs)
  - \* Universal Sentence Encoder with a minimum angular similarity of = 0.5.
  - \* Word Swapping Technique: Greedy Word Swap with Word Importance Ranking with word importance ranking conducted using input deletion method.

#### 6.4.2 TextBugger

- Word Importance Selection

- Max allowable replacement candidate generation for synonyms: 40.
- Transformation Embedding Mechanism: Counterfitted Glove Embeddings ([Mrkšić et al., 2016a](#))
- Allowable Swap Mechanisms: Character Insertion, Character Deletion, Adjacent Character Swap, Homoglyph Swap.

- Word Replacement:

- Pre-transformation constraints:
  - \* RepeatModification: A constraint disallowing the modification of words which have already been modified
  - \* StopwordModification: A constraint disallowing the modification of stopwords
- Constraints:
  - \* Universal Sentence Encoder with a minimum angular similarity of = 0.84
  - \* Word Swapping Technique: Greedy Word Swap with Word Importance Ranking with word importance ranking conducted using input deletion method.

#### 6.4.3 PWWS

- Word Importance Selection

- Max allowable replacement candidate generation for synonyms: 40.
- Transformation Embedding Mechanism: Word Swap by swapping synonyms in WordNet ([Miller, 1998](#))
- Allowable Swap Mechanisms: Character Insertion, Character Deletion, Adjacent Character Swap, Homoglyph Swap.

- Word Replacement:

- Pre-transformation constraints:
  - \* RepeatModification: A constraint disallowing the modification of words which have already been modified
  - \* StopwordModification: A constraint disallowing the modification of stopwords
- Constraints:
  - \* Word Swapping Technique: Greedy Word Swap with Word Importance Ranking with word importance ranking conducted using weighted saliency method.

#### 6.4.4 Kuleshov

- Word Importance Selection

- Max allowable replacement candidate generation for synonyms: 15.
- Transformation Embedding Mechanism: Counterfitted Glove Embeddings ([Mrkšić et al., 2016a](#))

- Word Replacement:

- Pre-transformation constraints:
  - \* RepeatModification: A constraint disallowing the modification of words which have already been modified
  - \* StopwordModification: A constraint disallowing the modification of stopwords
- Constraints:
  - \* Max words perturbed = 50
  - \* Maximum thought vector Euclidean distance = 0.2
  - \* Maximum language model log-probability difference = 2
  - \* Word Swapping Technique: Greedy Word Search.

### 6.4.5 DeepWordBug

- Word Importance Selection
  - Max allowable replacement candidate generation for synonyms: 40
  - Embedding Transformation Mechanism: Counterfitted Glove Embeddings (Mrkšić et al., 2016a)
  - Allowable Swap Mechanisms: Character Insertion, Character Deletion, Adjacent Character Swap, Random Character Substitution.
- Word Replacement:
  - Pre-transformation constraints:
    - \* RepeatModification: A constraint disallowing the modification of words which have already been modified
    - \* StopwordModification: A constraint disallowing the modification of stopwords
  - Constraints:
    - \* Maximum Levenshtien Edit Distance= 30.
    - \* Word Swapping Technique: Greedy Word Swap with Word Importance Ranking with word importance ranking conducted using input deletion method.

### 6.4.6 BAE

- Word Importance Selection
  - Max allowable replacement candidate generation for synonyms: 40
  - Transformation Embedding Mechanism: Transformer AutoTokenizer and word replacement using Masked Language Modelling. (Mrkšić et al., 2016a)
- Word Replacement:
  - Pre-transformation constraints:
    - \* RepeatModification: A constraint disallowing the modification of words which have already been modified
    - \* StopwordModification: A constraint disallowing the modification of stopwords
  - Constraints:
    - \* Part of Speech : Only replace words with the same part of speech (or nouns with verbs)

- \* Universal Sentence Encoder with a minimum angular similarity = 0.93.
- \* Word Swapping Technique: Greedy Word Swap with Word Importance Ranking with word importance ranking conducted using input deletion method.

## 6.5 Examples of NMT-TextAttack

1. **Original** : drawing on an irresistible , languid romanticism , byler reveals the ways in which a sultry evening or a beer-fueled afternoon in the sun can inspire even the most retiring heart to venture forth . (**Sentiment: Positive**)

**Adversarial (TextFooler)**: drawing on an **gar-gantuan** , **lolling melodrama** , byler **betrays** the ways in which a sultry evening or a beer-fueled afternoon in the sun can inspire even the most retiring heart to venture forth . (**Sentiment: Negative**)

**Adversarial (TextFooler+NMT-Text-Attack)**: drawing on an **inexorable**, **crooning melodrama** byler reveals the ways in which a sultry evening or a beer-fueled afternoon in the sun can inspire even the most retiring heart to venture forth. (**Sentiment: Negative**)

**Back-Translated (TextFooler)**: drawing on a giant melodrama, melodrama lolling, Byler betrays the ways in which a sensual afternoon or an afternoon of beer fed in the sun can inspire even the most outgoing heart to venture forward. (**Sentiment: Positive**)

**Back-Translated (TextFooler+NMT-Text-Attack)**: drawing on a melodrama byler **inexorable** betrays the ways in which a sensual afternoon or an afternoon of beer fed in the sun can inspire even the most outgoing heart to venture forward (**Sentiment: Negative**)

2. **Original** : Exceptionally well acted by Diane Lane and Richard Gere . (**Sentiment: Positive**)

**Adversarial (TextFooler)**: Exceptionally **op-portune** acted by Diane Lane and Richard Gere .(**Sentiment: Negative**)

**Adversarial (TextFooler+NMT-Text-Attack)**: Exceptionally **better** acted by Diane Lane and Richard Gere (**Sentiment: Negative**)

**Back-Translated (TextFooler)**: exceptionally timely performed by Diane Lane and Richard Gere. (**Sentiment: Positive**)

**Back-Translated (TextFooler+NMT-Text-Attack)**: exceptionally better performed by Diane Lane and Richard Gere (**Sentiment: Negative**)

**3.Original** : this kind of hands-on storytelling is ultimately what makes shanghai ghetto move beyond a good , dry , reliable textbook and what allows it to rank with its worthy predecessors . (**Sentiment: Positive**)

**Adversarial (PWWS)**: this **tolerant** of hands-on storytelling is ultimately what **piss** shanghai ghetto move beyond a good , dry , reliable textbook and what allows it to **gross** with its worthy predecessors (**Sentiment: Negative**)

**Adversarial (PWWS+NMT-TextAttack)**:this **tolerant** of hands-on storytelling is ultimately what makes shanghai ghetto move beyond a good , dry , reliable textbook and what allows it to **place** with its worthy predecessors . (**Sentiment: Negative**)

**Back-Translated (PWWS)**: This tolerant of practical narration is ultimately what **pis** shanghai ghetto move beyond a good, dry, reliable textbook and what allows rough with its worthy predecessors. (**Sentiment: Positive**)

**Back-Translated (PWWS+NMT-Text-Attack)**: this tolerant of narration is ultimately what builds the shanghai ghetto to move beyond a good reliable dry text book and what allows it to grossly with its worthy predecessors. (**Sentiment: Negative**)

**4.Original** : I went there today! I have an awful experience. They lady that cut my hair was nice but she wanted to leave early so she made a disaster in my head! (**Sentiment: Positive**)

**Adversarial (PWWS)**: I went there today! I have an **awesome** experience. They lady that cut my hair was nice but she wanted to leave early so she made a disaster in my head!(**Sentiment: Negative**)

**Adversarial (PWWS+NMT-TextAttack)**:I went there today! I have an **direful** experience! They lady that cut my hair was nice but she wanted to leave early so she made a disaster in my head (**Sentiment: Negative**)

**Back-Translated (PWWS)**: I went there today. I have a amazing experience. The lady who cut my hair was nice, but she wanted to leave early, so she made a mess of my head. (**Sentiment: Positive**)

**Back-Translated (PWWS+NMT-Text-Attack)**: I went there today. I have a terrible experience. The lady who cut my hair was nice, but she wanted to leave early, so she made a mess of my head.(**Sentiment: Negative**)

**5.Original** : I fell in love with this place as soon as we pulled up and saw the lights strung up and

oldies coming from the speakers! I tried the banana cream pie hard ice cream, their scoops are very generous!! My bf got the peach cobbler hard ice cream and that was to die for! We got 4 servings of ice cream for \$10, which nowadays is a steal IMO! :) I'll definitely be heading back with my coworkers this week! (**Sentiment: Positive**)

**Adversarial (TextBugger)**: I **declined** in **love** with this place as shortly as we **pulled** up and **saw** the headlights stung up and oldies coming from the speakers! I tried the **ban ana cream pe** hard ice cream, their scoops are very generous!! My bf got the peach cobbler hard ice cream and that was to die for! We got 4 servings of ice cream for \$10, which nowadays is a steal IMO! :) I'll **definitely** be heading back with my coworkers this **w eek**!(**Sentiment: Negative**)

**Adversarial (TextBugger+NMT-TextAttack)**:I fell in **love** with this place as soon as we pulled up and saw the lights strung up and oldies coming from the speakers! I tried the banana cream pie hard ice cream, their scoops are very generous!! My bf got the peach cobbler hard ice cream and that was to die for! We got 4 servings of ice cream for \$10, which existent is a theft IMO! :) I'll **doubtless** be heading back with my coworkers this week! (**Sentiment: Negative**)

**Back-Translated (TextBugger)**: I decided in love with this place as soon as we got up and climbed the chopped headlights and the old ones coming from the speakers! I've had the hard ice cream of ban ana cream, its spoonfuls are very generous! My friend got the hard iced peach pie and it was to die! We have 4 servings of ice cream for \$10, which today is an OMI robbery! :) I'll definitely be coming back with my coworkers this week! (**Sentiment: Positive**)

**Back-Translated (TextBugger+NMT-Text-Attack)**: I fell in love with this place as soon as we stopped and saw the stiff, old lights coming from the loudspeakers! I have tasted the hard frozen banana cream cake, its spoonfuls are very generous!! My bf got the hard iced peach pie and he was going to die for it! We have 4 servings of ice cream for \$10, which exists is an OMI robbery! :) I will definitely return with my co-workers this week!(**Sentiment: Negative**)



## 6.6 Walkthrough of TextFooler+NMT-Text-Attack

This section is concerned with providing an intuitive overview of the working of the attack agnostic NMT-Text-Attack algorithm with TextFooler. For ease of understanding, we use only one language for translation: Spanish. The algorithm, as shown before, is divided into sections. The first section, as shown in Figure 3, is the Word Importance Ranking section. Here, as per TextFooler's prescribed process, each word is replaced from the sentence and it's importance is evaluated by the change in classification score of the sentence before and after replacement.

On calculating the importance ranking score, we move to the second section, as shown in Figure 4. Here, we find synonyms for each word from the counterfitted GloVe word embeddings. These words are appended into the sentences replacing the original word, and passed to the NMT-Text-Attack Module. Here, the sentence undergoes round-trip translation to assess whether the inclusion of the word maintains robustness of the original attack under translation. We then collect the candidate sentences, and pass them through the final constraint requirement list, local to TextFooler. This includes checking whether the replaced word maintains the original word's POS tag, and then ranks them based on highest similarity score through USE embeddings and cosine similarity. Finally, we receive the adversarial example robust to round-trip translation.

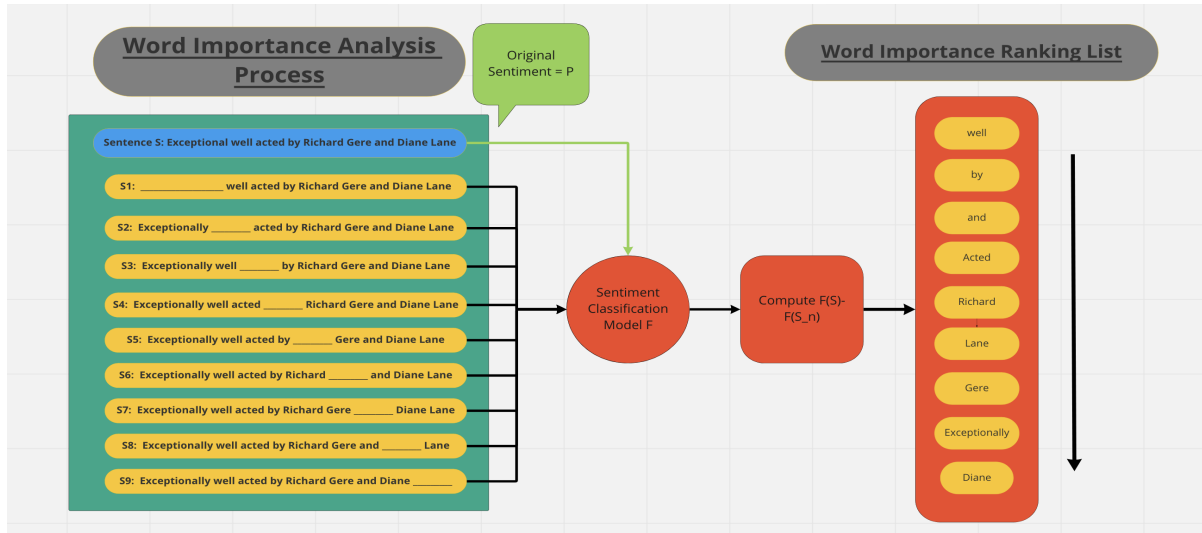


Figure 3: Word Importance Ranking Process

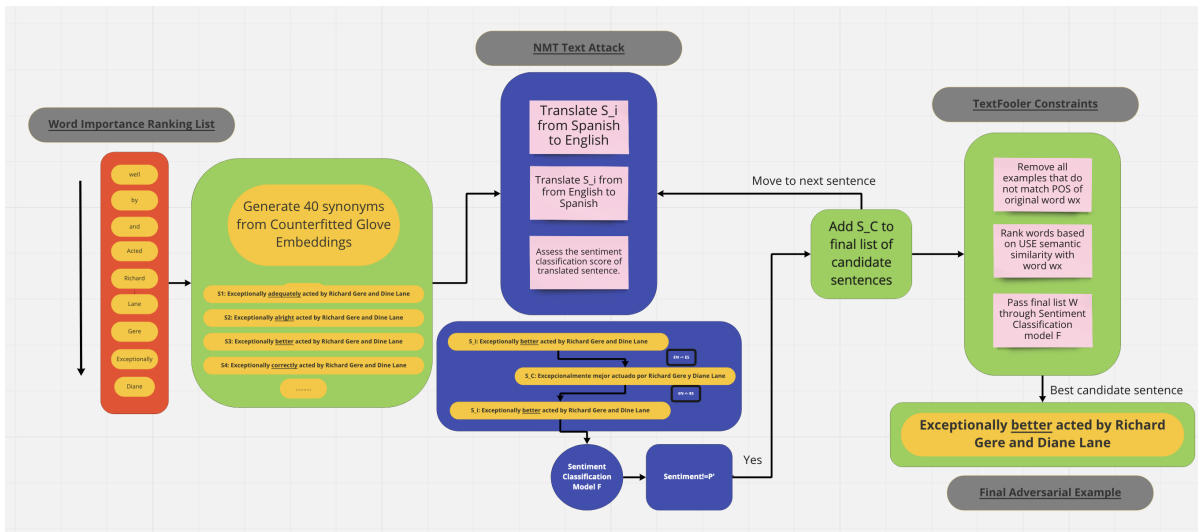


Figure 4: Word Replacement Process

## References

- Alfred V. Aho and Jeffrey D. Ullman. 1972. *The Theory of Parsing, Translation and Compiling*, volume 1. Prentice-Hall, Englewood Cliffs, NJ.
- Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018. [Generating natural language adversarial examples](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2890–2896, Brussels, Belgium. Association for Computational Linguistics.
- American Psychological Association. 1983. *Publications Manual*. American Psychological Association, Washington, DC.
- Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853.
- Galen Andrew and Jianfeng Gao. 2007. Scalable training of L1-regularized log-linear models. In *Proceedings of the 24th International Conference on Machine Learning*, pages 33–40.
- 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*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Jialun Cao, Meiziniu Li, Yeting Li, Ming Wen, and S. C. Cheung. 2020. Semmt: A semantic-based testing approach for machine translation systems. *ArXiv*, abs/2012.01815.
- Daniel Cer, Yinfei Yang, Sheng yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. [Universal sentence encoder](#).
- Sin-wai Chan. 2006. *A Dictionary of Translation Technology*. The Chinese University of Hong Kong Press.
- Ashok K. Chandra, Dexter C. Kozen, and Larry J. Stockmeyer. 1981. [Alternation](#). *Journal of the Association for Computing Machinery*, 28(1):114–133.
- Aochuan Chen, Yuguang Yao, Pin-Yu Chen, Yihua Zhang, and Sijia Liu. 2023. Understanding and improving visual prompting: A label-mapping perspective. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Pin-Yu Chen and Payel Das. 2023. Ai maintenance: A robustness perspective. *IEEE Communications Magazine*.
- Pin-Yu Chen and Cho-Jui Hsieh. 2023. *Adversarial Robustness for Machine Learning*. Elsevier.
- Pin-Yu Chen and Sijia Liu. 2023. Holistic adversarial robustness of deep learning models. *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Minhao Cheng, Jinfeng Yi, Pin-Yu Chen, Huan Zhang, and Cho-Jui Hsieh. 2020. [Seq2sick: Evaluating the robustness of sequence-to-sequence models with adversarial examples](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 3601–3608. AAAI Press.
- Yong Cheng, Lu Jiang, and Wolfgang Macherey. 2019. [Robust neural machine translation with doubly adversarial inputs](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4324–4333, Florence, Italy. Association for Computational Linguistics.
- Zewei Chu, Mingda Chen, Jing Chen, Miaosen Wang, Kevin Gimpel, Manaal Faruqui, and Xiance Si. 2020. [How to ask better questions? A large-scale multi-domain dataset for rewriting ill-formed questions](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7586–7593. AAAI Press.
- Nathan E. Crone, A. J. Power, and John Weldon. 2021. Quality estimation using round-trip translation with sentence embeddings. *ArXiv*, abs/2111.00554.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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 (Volume 2: Short Papers)*, pages 31–36, Melbourne, Australia. 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](#).
- Siddhant Garg and Goutham Ramakrishnan. 2020. [BAE: BERT-based adversarial examples for text classification](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6174–6181, Online. Association for Computational Linguistics.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. [Explaining and harnessing adversarial examples](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Zilu Guo, Zhongqiang Huang, Kenny Q. Zhu, Guandan Chen, Kaibo Zhang, Boxing Chen, and Fei Huang. 2021. Automatically paraphrasing via sentence reconstruction and round-trip translation. In *IJCAI*.
- Dan Gusfield. 1997. *Algorithms on Strings, Trees and Sequences*. Cambridge University Press, Cambridge, UK.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural Comput.*, 9(8):1735–1780.
- 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)*, pages 4129–4142, Hong Kong, China. Association for Computational Linguistics.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. [Is BERT really robust? A strong baseline for natural language attack on text classification and entailment](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8018–8025. AAAI Press.
- Volodymyr Kuleshov, Shantanu Thakoor, Tingfung Lau, and Stefano Ermon. 2018. Adversarial examples for natural language classification problems.
- Qi Lei, Lingfei Wu, Pin-Yu Chen, Alexandros G Dimakis, Inderjit S Dhillon, and Michael Witbrock. 2019. Discrete adversarial attacks and submodular optimization with applications to text classification. *SysML*.
- Dianqi Li, Yizhe Zhang, Hao Peng, Liqun Chen, Chris Brockett, Ming-Ting Sun, and Bill Dolan. 2021. [Contextualized perturbation for textual adversarial attack](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5053–5069, Online. Association for Computational Linguistics.
- Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2019. Textbugger: Generating adversarial text against real-world applications. *ArXiv*, abs/1812.05271.
- Jiwei Li, Will Monroe, and Dan Jurafsky. 2017. [Understanding neural networks through representation erasure](#).
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. [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, Online. 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.

- Jared Lichtarge, Chris Alberti, Shankar Kumar, Noam Shazeer, Niki Parmar, and Simon Tong. 2019. [Corpora generation for grammatical error correction](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3291–3301, Minneapolis, Minnesota. Association for Computational Linguistics.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. [Learning word vectors for sentiment analysis](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.
- Nitin Madnani, Joel Tetreault, and Martin Chodorow. 2012. [Exploring grammatical error correction with not-so-crummy machine translation](#). In *Proceedings of the Seventh Workshop on Building Educational Applications Using NLP*, pages 44–53, Montréal, Canada. Association for Computational Linguistics.
- George A Miller. 1998. *WordNet: An electronic lexical database*. MIT press.
- Jihyung Moon, Hyunchang Cho, and Eunjeong L. Park. 2020. [Revisiting round-trip translation for quality estimation](#). In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 91–104, Lisboa, Portugal. European Association for Machine Translation.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. 2017. [Universal adversarial perturbations](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 86–94. IEEE Computer Society.
- John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. [TextAttack: A framework for adversarial attacks, data augmentation, and adversarial training in NLP](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 119–126, Online. 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. 2016a. [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.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Blaise Thomson, Milica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, David Vandyke, Tsung-Hsien Wen, and Steve Young. 2016b. [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.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Blaise Thomson, Milica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, David Vandyke, Tsung-Hsien Wen, and Steve Young. 2016c. [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.
- Bo Pang and Lillian Lee. 2005. [Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales](#). In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 115–124, Ann Arbor, Michigan. Association for Computational Linguistics.
- Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. 2017. [Practical black-box attacks against machine learning](#).
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015. [Yara parser: A fast and accurate dependency parser](#). *Computing Research Repository*, arXiv:1503.06733. Version 2.
- Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. [Generating natural language adversarial examples through probability weighted word saliency](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1085–1097, Florence, Italy. Association for Computational Linguistics.



- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. [Semantically equivalent adversarial rules for debugging NLP models](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 856–865, Melbourne, Australia. Association for Computational Linguistics.
- Suranjana Samanta and Sameep Mehta. 2017. [Towards crafting text adversarial samples](#).
- Tomohiro Shigenobu. 2007. Evaluation and usability of back translation for intercultural communication. In *Usability and Internationalization. Global and Local User Interfaces*, pages 259–265, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Harold Somers. 2005. [Round-trip translation: What is it good for?](#) In *Proceedings of the Australasian Language Technology Workshop 2005*, pages 127–133, Sydney, Australia.
- Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. [Multilingual translation with extensible multilingual pretraining and finetuning](#).
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Huggingface’s transformers: State-of-the-art natural language processing](#).
- 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*, pages 6066–6080, Online. Association for Computational Linguistics.
- Zhiyuan Zeng and Deyi Xiong. 2021a. [An empirical study on adversarial attack on NMT: Languages and positions matter](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 454–460, Online. Association for Computational Linguistics.
- Zhiyuan Zeng and Deyi Xiong. 2021b. [An empirical study on adversarial attack on NMT: Languages and positions matter](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 454–460, Online. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 649–657.
- Xinze Zhang, Junzhe Zhang, Zhenhua Chen, and Kun He. 2021. [Crafting adversarial examples for neural machine translation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1967–1977, Online. Association for Computational Linguistics.
- 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.