# BeamAttack: Generating High-quality Textual Adversarial Examples through Beam Search and Mixed Semantic Spaces

Hai Zhu[1,3,*]([⊠]), Qinyang Zhao[2], and Yuren Wu[3]

[1] University of Science and Technology of China, Hefei, China
SA21218029@mail.ustc.edu.cn
[2] Xidian University, Xi'an, China
21151213588@stu.xidian.edu.cn
[3] Ping An Technology (Shenzhen) Co., Ltd., Shenzhen, China
wuyuren134@pingan.com.cn

**Abstract.** Natural language processing models based on neural networks are vulnerable to adversarial examples. These adversarial examples are imperceptible to human readers but can mislead models to make the wrong predictions. In a black-box setting, attacker can fool the model without knowing model's parameters and architecture. Previous works on word-level attacks widely use single semantic space and greedy search as a search strategy. However, these methods fail to balance the attack success rate, quality of adversarial examples and time consumption. In this paper, we propose **BeamAttack**, a textual attack algorithm that makes use of mixed semantic spaces and improved beam search to craft high-quality adversarial examples. Extensive experiments demonstrate that BeamAttack can improve attack success rate while saving numerous queries and time, e.g., improving at most 7% attack success rate than greedy search when attacking the examples from MR dataset. Compared with heuristic search, BeamAttack can save at most 85% model queries and achieve a competitive attack success rate. The adversarial examples crafted by BeamAttack are highly transferable and can effectively improve model's robustness during adversarial training. Code is available at https://github.com/zhuhai-ustc/beamattack/tree/master

**Keywords:** Adversarial Examples · Robustness · Natural Language Processing

## 1 Introduction

In recent years, neural networks have achieved great success in the natural language processing field while being vulnerable to adversarial examples. These adversarial examples are original inputs altered by some tiny perturbations [9,

---

* This work was done when the author was at Ping An Technology (Shenzhen) Co., Ltd.

20]. It is worth noting that perturbations are imperceptible to humans but can mislead the model decision. Therefore, it is essential to explore adversarial examples since our goal is to improve the reliability and robustness of the model, especially on some security-critical applications, such as toxic text detection and public opinion analysis [25]. Compared to image and speech attacks [22, 2], it is more challenging in crafting textual adversarial examples due to the discrete of natural language. In addition, there are some grammar constraints in the textual adversarial examples:(1)the crafted examples should keep the same meaning as the original texts,(2)generated examples should look natural and grammatical. However, previous works barely conform to all constraints, or satisfy the above constraints at the cost of reducing the attack success rate.

Conventional word-level attack algorithms can be roughly divided into three steps: (1) calculating word importance score according to the changes of class label probabilities after replacing this word, (2) searching synonyms for each origin word, (3) selecting the substitution that reduces the class label probabilities most and replacing origin word with it until model predicts wrong. The problem is that previous works only use a single semantic space to search synonyms, which limits the diversity of substitutions and cut down the search space. In addition, most prior works introduce greedy search to select the best substitution with the maximum change of class label probabilities [9, 11]. Greedy search limits the search space and sometimes leads to local optimal solution and word over-substitution. Therefore, some works [26, 20] introduce heuristic search to improve attack success rate, at the cost of time-consuming and numerous model queries. In generally, previous works fail to balance the attack success rate, quality of adversarial examples and time consumption.

In this paper, we propose BeamAttack, a textual attack algorithm based on mixed semantic spaces and beam search. Specially, we search substitutions from word embedding space and BERT respectively, and filter out the bad synonyms to improve semantic similarity of adversarial examples, then improve beam search to craft adversarial examples, which greatly expands the search space by controlling beam size. Therefore, it is capable of escaping from local optima within acceptable number of model queries. Furthermore, we evaluate BeamAttack by attacking various neural networks on five datasets. Experiments show that it outperforms other baselines in attack success rate and semantic similarity while saving numerous model queries. Our main contributions are summarized as follows:

- We propose the mixed semantic spaces, making full use of word embedding space and BERT simultaneously to expand the diversity of substitutions and generating high-qualify adversarial examples.
- We propose BeamAttack, a black-box attack algorithm which improves beam search to expand search space and reduce the redundancy word substitution.
- Experiments show that BeamAttack achieves the trade-off results compared with previous works. In addition, adversarial examples crafted by BeamAttack with high semantic similarity, low perturbation, and good transferability.

## 2    Related Work

We divide the existing textual attack algorithms into char-level, word-level and sentence-level attacks based on granularity. Char-level attacks generate adversarial examples by inserting, swapping or removing characters(such as 'attack' → 'atttack') [12, 6], which can be easily rectified by word spelling machine. Sentence-level attacks insert some perturbed sentences into the origin paragraph to confuse models [3]. Nevertheless, these adversarial examples contain many lexical errors.

In order to generate high-quality adversarial examples, word-level attacks have gradually become a prevalent approach. Word-level attacks substitute the origin words with synonyms(such as 'like' → 'love'). Traditional strategies search synonyms from word embedding space. For example, some works [14, 9, 23] calculate the word saliency and greedily substitute words with synonyms derived from WordNet [16], or utilizing word importance score and replace words with synonyms from counter-fitting word vectors [18]. Recently, researcher [13, 7, 11] search synonyms from pre-trained language models (e.g. BERT, RoBERTa). The pre-trained language models are trained on massive text data, and predict the masked words. Therefore, it has the ability to predict contextual-aware words.

Above attack algorithms adopt the greedy search, which limits the search space and leads to local optimal solution. Minor work have explored the heuristic search, such as genetic algorithm [20],particle swarm optimization [26]. However, heuristic search is very time-consuming and requires a lot of model queries. Therefore, we propose BeamAttack, searching synonyms from word embedding space and BERT simultaneously, and fine-tuning beam search to expand search space and reduce word-over substitution.

## 3    Beam Search Adversarial Attack

BeamAttack is divided into three steps. There are word importance calculation, mixed semantic spaces and improved beam search. The overview of BeamAttack is shown in the Figure 1. Before delving into details, we present the attack settings and problem formulation.

### 3.1    Black-box Untargeted Attack

The BeamAttack belongs to black-box attacks, it has nothing about model's architecture, parameters and gradients, only class label probabilities are accessible. Given a sentence of $n$ words $\mathcal{X} = [x_1, x_2, \cdots, x_n]$ and label set $\mathcal{Y}$, a well-trained model can classify sentence correctly:

$$\underset{y_i \in \mathcal{Y}}{\operatorname{argmax}} P(y_i|\mathcal{X}) = y_{true} \tag{1}$$

The adversarial example $\mathcal{X}' = [x_1', x_2', \cdots, x_n']$ is crafted to make model predict wrong. In addition, there are some constraints on the word substitution
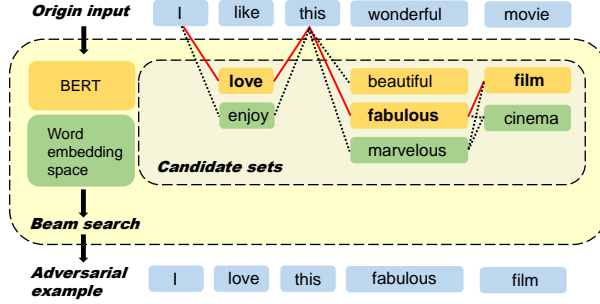
**Fig. 1.** The overview of BeamAttack. Candidate sets are substitutions generated from BERT and word embedding space. Black lines are beam search paths, wherein red lines are the optimal search path.

rate(WSR) and semantic similarity(SIM) of the adversarial example. $\mathcal{X}'$ should be close to $\mathcal{X}$ and a human reader hardly differentiate the modifications. The mathematical expression is as follows:

$$\underset{y_i \in \mathcal{Y}}{\operatorname{argmax}} P(y_i | \mathcal{X}') \neq y_{true}$$
$$\text{s.t. } \mathrm{SIM}(\mathcal{X}', \mathcal{X}) > \mathcal{L}; \mathrm{WSR}(\mathcal{X}', \mathcal{X}) < \sigma \tag{2}$$

### 3.2   Word Importance Calculation

Given a sentence of $n$ words $\mathcal{X} = [x_1, x_2, \cdots, x_n]$, only some important words will affect the prediction results of the model $\mathcal{F}$. In order to measure the importance of $x_i$, we follow the calculation proposed in TextFooler [9]. We replace $x_i$ with '[oov]'[4] to form $\mathcal{X}/\{x_i\} = [x_1, \cdots, x_{i-1}, [oov], x_{i+1}, \cdots, x_n]$, then word importance of $x_i$ is calculated as follows:

- The predicted label remains the same after replace, i.e., $\mathcal{F}(\mathcal{X}) = \mathcal{F}(\mathcal{X}/\{x_i\}) = y_{true}$,

$$I(x_i) = \mathcal{F}_{y_{true}}(\mathcal{X}) - \mathcal{F}_{y_{true}}(\mathcal{X}/\{x_i\}) \tag{3}$$

- The predicted label is changed after replace, i.e., $\mathcal{F}(\mathcal{X}) = y_{true} \neq y_{other} = \mathcal{F}(\mathcal{X}/\{x_i\})$,

$$I(x_i) = \mathcal{F}_{y_{true}}(\mathcal{X}) - \mathcal{F}_{y_{true}}(\mathcal{X}/\{x_i\}) + \mathcal{F}_{y_{other}}(\mathcal{X}/\{x_i\}) - \mathcal{F}_{y_{other}}(\mathcal{X}) \tag{4}$$

where $\mathcal{F}_y(\mathcal{X})$ represents the predicted class label probability of $\mathcal{X}$ by $\mathcal{F}$ on label $y$. In order to improve the readability and fluency of the adversarial examples, we will filter out stopwords by NLTK[5] after calculating the word importance.

---

[4] the word out-of-vocabulary.
[5] https://www.nltk.org/

### 3.3   Mixed Semantic Spaces

After ranking the words by their importance score, we need to search synonyms, which is a candidate words set $\mathcal{C}(x_i)$ for each word $x_i$. A proper replacement word should (i) have similar semantic meaning with original input, (ii) avoid some obvious grammar errors, (iii) and confuse model $\mathcal{F}$ to predict the wrong label. There are two different semantic spaces to search synonyms, word embedding spaces and pre-trained language models.

– The former searches for synonyms from word embedding spaces, such as WordNet space [16], HowNet space [5] and Counter-fitting word vectors [18]. Word embedding spaces can quickly generate synonyms with the same meaning as origin word.
– The later searches for synonyms through pre-trained language models(such as BERT). Given a sentence of $n$ words $\mathcal{X} = [x_1, x_2, \cdots, x_n]$, we replace each word $x_i$ with '[MASK]', and get candidate words set $\mathcal{C}(x_i)$ predicted by BERT. Pre-trained language models produce fluent and contextual-aware adversarial examples.

We combine word embedding space and BERT to make full use of the advantage of different semantic spaces. In detail, for each word $x_i$, we respectively select top $N$ synonyms from word embedding space and BERT to form a candidate words set $\mathcal{C}(x_i)$. To generate high-qualify adversarial examples, we filter out the candidate words set that has different part-of-speech(POS)[6] synonyms with $x_i$. In addition, for each $c \in \mathcal{C}(x_i)$, we substitute it for $x_i$ to generate adversarial example $\mathcal{X}' = [x_1, \cdots, x_{i-1}, c, x_{i+1}, \cdots, x_n]$, then we measure semantic similarity between $\mathcal{X}$ and adversarial example $\mathcal{X}'$ by universal sentence encoder(USE)[7], which encodes original input $\mathcal{X}$ and adversarial example $\mathcal{X}'$ as dense vectors and use cosine similarity as a approximation of semantic similarity. Only synonyms whose similarity is higher than threshold $L$ will be retained in the candidate words set $\mathcal{C}(x_i)$.

### 3.4   Improved Beam Search

After filtering out the candidate words set $\mathcal{C}(x_i)$, the construction of adversarial examples is a combinatorial optimization problem as expected in Eq.2. Previous works use the greedy search since it solely selects the token that maximizes the probability difference, which leads to local optima and word-over substitution.

To tackle this, we improve beam search to give consideration to both attack success rate and algorithm efficiency. Beam search has a hyper-parameter called beam size $\mathcal{K}$. Naive beam search only selects top $\mathcal{K}$ adversarial examples from the current iteration results. In the improved beam search, we merge the output of the last iteration to the current iteration and select top $\mathcal{K}$ adversarial examples as the input of the next iteration jointly. In detail, for each word $x_i$ in the original

---

[6] https://spacy.io/api/tagger
[7] https://tfhub.dev/google/ universal-sentence-encoder

text, we replace $x_i$ with the substitution from candidate words set $\mathcal{C}(x_i)$ to generate adversarial examples $\mathcal{X}'$ and calculate the probability differences. The top $\mathcal{K}$ adversarial examples $\mathcal{X}'$ with the maximum probability difference(including the last iteration of top $\mathcal{K}$ adversarial examples) are selected as the input of the next iteration until the attack succeeds or all origin words are iterated. It is worth noting that greedy search is a special case of $\mathcal{K} = 1$. The details of BeamAttack are shown in algorithm 1.

---

**Algorithm 1** BeamAttack Adversarial Algorithm

---

**Input**:Original text $\mathcal{X}$, target model $\mathcal{F}$, semantic similarity threshold $\mathcal{L} = 0.5$ and beam size $\mathcal{K} = 10$, number of words in original text $n$
**Output**:Adversarial example $\mathcal{X}_{\text{adv}}$.

1: $\mathcal{X}_{\text{adv}} \leftarrow \mathcal{X}$
2: $set(\mathcal{X}_{\text{adv}}) \leftarrow \mathcal{X}_{\text{adv}}$
3: **for** each word $x_i$ in $\mathcal{X}$ **do**
4:     Compute the importance score $I(x_i)$ via Eq.3 and 4.
5: **end for**
6: Sort the words with importance score $I(x_i)$
7: **for** $i = 1$ to $n$ **do**
8:     Replace the $x_i$ with [MASK]
9:     Generate the candidate set $\mathcal{C}(x_i)$ from BERT and Word Embedding Space
10:     $\mathcal{C}(x_i) \leftarrow \text{POSFilter}(\mathcal{C}(x_i)) \cap \text{USEFilter}(\mathcal{C}(x_i))$
11: **end for**
12: **for** $\mathcal{X}_{\text{adv}}$ in $set(\mathcal{X}_{\text{adv}})$ **do**
13:     **for** $c_k$ in $\mathcal{C}(x_i)$ **do**
14:         $\mathcal{X}'_{\text{adv}} \leftarrow$ Replace $x_i$ with $c_k$ in $\mathcal{X}_{\text{adv}}$
15:         Add $\mathcal{X}'_{\text{adv}}$ to the $set(\mathcal{X}_{\text{adv}})$
16:     **end for**
17:     **for** $\mathcal{X}'_{\text{adv}}$ in $set(\mathcal{X}_{\text{adv}})$ **do**
18:         **if** $\mathcal{F}(\mathcal{X}'_{\text{adv}}) \neq y_{true}$ **then**
19:             **return** $\mathcal{X}'_{\text{adv}}$ with highest semantic similarity
20:         **end if**
21:     **end for**
22:     $set(\mathcal{X}_{\text{adv}}) \leftarrow$ Select top $\mathcal{K}$ adversarial examples in $set(\mathcal{X}_{\text{adv}})$
23:     $i \leftarrow i + 1$
24:     **if** $i > n$ **then**
25:         **break**
26:     **end if**
27: **end for**
28: **return** adversarial examples $\mathcal{X}_{\text{adv}}$

---

## 4   Experiments

**Tasks, Datasets and Models.** To evaluate the effectiveness of BeamAttack, we conduct experiments on two NLP tasks, including text classification and

text inference. In particular, the experiments cover various datasets, such as MR [19],IMDB [15],SST-2 [21], SNLI [1] and MultiNLI [24]. We train three neural networks as target models including CNN [10], LSTM [8] and BERT [4]. Model parameters are consistent with TextFooler's[9] setting.

**Baselines.** To quantitatively evaluate BeamAttack, we compare it with other black-box attack algorithms, including TextFooler(TF) [9], PWWS [20], BAE [7], Bert-Attack(BEAT) [13] and PSO [26], wherein TF,PWWS and PSO search synonyms from word embedding spaces, BAE and BEAT search synonyms from BERT. In addition, PSO belongs to heuristics search and other belong to greedy search. These baselines are implemented on the TextAttack framework [17].

**Automatic Evaluation Metrics.** We evaluate the attack performance by following metrics. Attack Success Rate(ASR) is defined as the proportion of successful adversarial examples to the total number of examples. Word Substitution Rate(WSR) is defined as the proportion of number of replacement words to number of origin words. Semantic Similarity(SIM) is measured by Universal Sentence Encoder(USE). Query Num(Query) is the number of model queries during adversarial attack. The ASR evaluates how successful the attack is. The WSR and semantic similarity together evaluate how semantically similar the original texts and adversarial examples are. Query num can reveal the efficiency of the attack algorithm.

**Implementation Details.** In our experiments, we carry out all experiments on NVIDIA Tesla P100 16G GPU. We set the beam size $\mathcal{K} = 10$, number of each candidate set $N = 50$, semantic similarity threshold $\mathcal{L} = 0.5$, we take the average value of 1000 examples as the final experimental result.

## 4.1 Experimental Results

The experiment results are listed in Table 1. It is worth noting that BeamAttack achieves higher ASR than baselines on almost all scenarios. BeamAttack also reduces the WSR on some datasets(MR,IMDB and SST-2). We attribute this superiority to the fine-tuned beam search, as this is the major improvement of our algorithm compared with greedy search. BeamAttack has chance to jump out of the local optimal solution and find out the adversarial examples with lower perturbation by expanding the search space. In terms of model robustness, BERT has better robustness than traditional classifiers(CNN and LSTM), since the attack success rate of attacking BERT is lower than other models.

**Semantic Similarity.** Except ASR and WSR, fluent and contextual-aware adversarial examples are also essential. Figure 2 plots the semantic similarity of adversarial examples generated by different attack algorithms. Clearly, BeamAttack achieves the highest semantic similarity than other attack algorithms.

**Qualitative Examples.** To more intuitively contrast the fluency of adversarial examples, we list some adversarial examples generated by different attack algorithms in Table 2. Compared with other methods, Beamattack not only ensures the semantic similarity between replacement words and original words but also successfully misleads the model with the minimum perturbation.

**Table 1.** The attack success rate and word substitution rate of different attack algorithms on five datasets. The "Origin ACC(%)" denotes the target model's test accuracy on the original inputs.

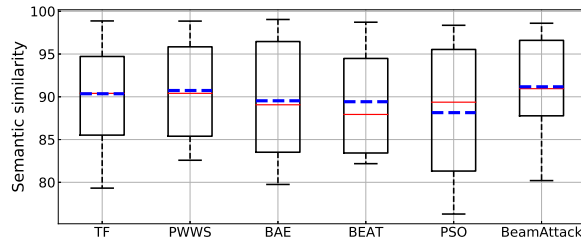| Datasets | Target Models | Origin ACC | Attack Success Rate(ASR(%)) | | | | | | Word Substitution Rate(WSR(%)) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | TF | PWWS | BAE | BEAT | PSO | BeamAttack | TF | PWWS | BAE | BEAT | PSO | BeamAttack |
| MR | CNN | 80.4 | 98.81 | 98.61 | 98.00 | 83.31 | 96.23 | **99.87** | 17.05 | 13.22 | 12.98 | 15.06 | 11.53 | **8.29** |
| | LSTM | 80.7 | 98.92 | 97.92 | 98.21 | 84.12 | 95.32 | **99.90** | 15.61 | 13.07 | 11.71 | 13.59 | 10.91 | **8.60** |
| | BERT | 90.4 | 90.54 | 81.53 | 90.61 | 88.36 | 92.47 | **97.88** | 20.91 | 14.67 | 14.44 | 15.32 | 11.93 | **9.70** |
| IMDB | CNN | 89.2 | **100** | **100** | **100** | 99.82 | **100** | **100** | 2.51 | 2.23 | 2.01 | 3.32 | 2.43 | **2.11** |
| | LSTM | 89.8 | 99.76 | 99.47 | **100** | 99.83 | **100** | **100** | 3.12 | 3.11 | **2.25** | 3.45 | 2.46 | 2.43 |
| | BERT | 90.9 | 88.83 | 86.55 | 83.96 | 88.68 | 89.93 | **91.6** | 3.81 | 5.02 | 7.69 | 5.66 | 4.32 | **1.65** |
| SST-2 | CNN | 82.5 | 92.37 | 98.23 | 95.45 | 86.44 | 96.69 | **99.88** | 17.09 | 13.10 | 12.53 | 15.40 | 11.47 | **8.46** |
| | LSTM | 84.6 | 93.21 | 98.48 | 96.23 | 86.43 | 96.42 | **100** | 17.55 | 13.53 | 12.83 | 15.31 | 11.45 | **8.76** |
| SNLI | BERT | 89.1 | 96.00 | 98.42 | 98.84 | 98.64 | 92.51 | **99.80** | 17.26 | 13.72 | **6.91** | 7.80 | 8.19 | 13.81 |
| MNLI | BERT | 85.1 | 90.44 | 94.33 | 99.23 | 92.00 | 83.43 | **99.50** | 13.93 | 10.12 | **5.45** | 5.64 | 6.65 | 10.81 |



**Fig. 2.** The semantic similarity between origin inputs and adversarial examples.

**Model Query.** The number of model queries measures the effectiveness of attack algorithm. Table 3 lists the model queries of various attack algorithms. Results show that although our BeamAttack needs more model queries than greedy search(such as TF), compared with the PSO attack algorithm, which adopts heuristic search, our algorithm obtains competitive results with extremely few model queries.

## 4.2   Ablation Study

**The effect of beam size $\mathcal{K}$.** To validate the effectiveness of beam size $\mathcal{K}$, we use BERT as the target model and test on MR dataset with different beam size $\mathcal{K}$. When $\mathcal{K} = 1$, beam search is equal to greedy search. As shown in Table 4, the attack success rate increases gradually with the grow of beam size $\mathcal{K}$.

**The effect of mixed semantic spaces.** Another major improvement of our BeamAttack is that substitutions are selected from mixed semantic spaces. As

**Table 2.** The adversarial example crafted by different attack algorithms on MR(BERT) dataset. Replacement words are represented in red.

| | |
|---|---|
| **Origin Text (Positive)** | The experience of the roles in the play makes us generates an enormous feeling of empathy for its characters. |
| **BAE (Negative)** | The experience of the roles in the play makes us generates an **excessive need** of empathy for its characters. |
| **PWWS (Negative)** | The experience of the roles in the play makes us **render** an enormous **smell** of empathy for its **eccentric**. |
| **TextFooler (Negative)** | The experience of the roles in the play makes us **leeds** an enormous **foreboding** of empathy for its **specs**. |
| **BeamAttack (Negative)** | The experience of the roles in the play makes us generates an enormous feeling of **pity** for its characters. |

**Table 3.** The average model queries of different attack algorithms on five datasets. Beam size $\mathcal{K} = 10$

| | MR | IMDB | SST-2 | SNLI | MNLI |
|---|---|---|---|---|---|
| TF | 113.8 | 536.7 | 146.2 | 54.1 | 68.9 |
| PWWS | 285.4 | 3286.5 | 5054.3 | 137.7 | 157.4 |
| BAE | 104.2 | 567.1 | 171.0 | 75.5 | 75.1 |
| BEAT | 207.9 | 585.0 | 245.6 | 93.6 | 119.2 |
| PSO | 5124.5 | 15564.3 | 3522.8 | 416.6 | 1124.8 |
| BeamAttack | 650.3 | 2135.8 | 584.3 | 126.0 | 174.0 |

shown in the Table 5, we study the impact of different semantic spaces on different metrics. Compared with single word embedding space or BERT, using both word embedding space and BERT to generate adversarial examples can obtain higher attack success rate, semantic similarity and lower word substitution rate.

### 4.3   Transferability

The transferability of adversarial examples reflects property that adversarial examples crafted by classifier $\mathcal{F}$ can also fool other unknown classifier $\mathcal{F}'$. We evaluate the transferability on MR dataset across CNN,LSTM and BERT. In detail, we use the adversarial examples crafted for attacking BERT on MR dataset to evaluate the transferability for CNN and LSTM models. As shown in the Figure 3, the adversarial examples generated by BeamAttack achieve the higher transferability than baselines.

### 4.4   Adversarial Training

Adversarial training is a prevalent technique to improve the model's robustness by adding some adversarial examples into train data. To validate this, we train the CNN model on the MR dataset and obtains 80.4% test accuracy. Then we

**Table 4.** The effect of beam size $\mathcal{K}$ on MR(BERT) dataset.

| Beam Size | ASR(%) | WSR(%) | Similarity(%) | Query |
|-----------|--------|--------|---------------|-------|
| $\mathcal{K} = 1$ | 89.0 | 15.5 | 81.3 | **101.3** |
| $\mathcal{K} = 2$ | 90.6 | 15.3 | 82.0 | 150.4 |
| $\mathcal{K} = 5$ | 91.6 | **15.1** | 82.8 | 312.6 |
| $\mathcal{K} = 7$ | 92.3 | **15.1** | 83.0 | 411.1 |
| $\mathcal{K} = 10$ | **92.6** | **15.1** | **83.1** | 516.2 |

**Table 5.** The effect of different semantic spaces on MR(BERT) dataset.

| semantic space | ASR(%) | WSR(%) | Similarity(%) | Query |
|----------------|--------|--------|---------------|-------|
| Embedding | 89.0 | 15.3 | 81.2 | 101.3 |
| BERT | 93.6 | 13.1 | 82.6 | **101.1** |
| Embedding+BERT | **95.3** | **11.7** | **84.9** | 140.3 |

randomly generate 1000 MR adversarial examples to its training data and re-train the CNN model. The result is shown in the Table 6, CNN model obtains 83.3% test accuracy, higher than origin test accuracy. Although there is no significant change in ASR, BeamAttack needs to replace more words and more model queries to attack successfully with WSR and model queries increasing. It indicates that adversarial training effectively improves the generalization and robustness of the model.

**Table 6.** The performance of CNN with(out) adversarial training on the MR dataset.

| | Origin ACC(%) | ASR(%) | WSR(%) | SIM(%) | Query |
|---|---------------|--------|--------|--------|-------|
| **Original** | 80.4 | 99.87 | 8.20 | 91.08 | 563.1 |
| **Adv.Training** | **83.3** | **99.75** | **8.67** | **90.82** | **606.4** |

## 5   Conclusion

In this paper, we propose an efficient adversarial textual attack algorithm BeamAttack. The BeamAttack makes full use of word embedding space and BERT to generate substitutions and fine-tune beam search to expand search spaces. Extensive experiments demonstrate BeamAttack balances the attack success rate, qualify of adversarial examples and time consumption. In addition, the adversarial examples crafted by BeamAttack are contextual-aware and improve models' robustness during adversarial training.
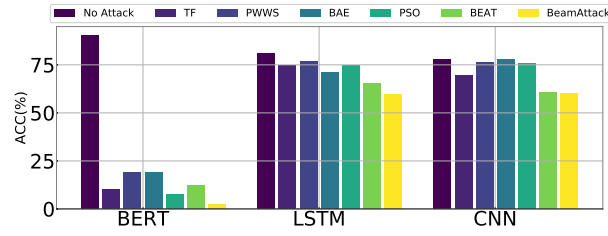
**Fig. 3.** Transfer attack on MR dataset. Lower accuracy indicates higher transferability (the lower the better).

# References

1. Bowman, S.R., Angeli, G., Potts, C., Manning, C.D.: large annotated corpus for learning natural language inference. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. pp. 632–642 (2015)
2. Carlini, N., Wagner, D.: Audio adversarial examples: Targeted attacks on speech-to-text. In: 2018 IEEE Security and Privacy Workshops. pp. 1–7 (2018)
3. Cheng, M., Yi, J., Chen, P.Y., Zhang, H., Hsieh, C.J.: Seq2sick: Evaluating the robustness of sequence-to-sequence models with adversarial examples. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 3601–3608 (2020)
4. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 4171–4186 (2019)
5. Dong, Z., Dong, Q., Hao, C.: Hownet and its computation of meaning. In: Proceedings of the 23rd International Conference on Computational Linguistics: Demonstrations. pp. 53–56 (2010)
6. Gao, J., Lanchantin, J., Soffa, M.L., Qi, Y.: Black-box generation of adversarial text sequences to evade deep learning classifiers. In: 2018 IEEE Security and Privacy Workshops. pp. 50–56 (2018)
7. Garg, S., Ramakrishnan, G.: BAE: BERT-based adversarial examples for text classification. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing. pp. 6174–6181 (2020)
8. Hochreiter, S., Schmidhuber, J.: Long short-term memory. In: Neural Computation. vol. 9, pp. 1735–1780 (1997)
9. Jin, D., Jin, Z., Zhou, J.T., Szolovits, P.: Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 8018–8025 (2020)
10. Kim, Y.: Convolutional neural networks for sentence classification. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing. pp. 1746–1751 (2014)
11. Li, D., Zhang, Y., Peng, H., Chen, L., Brockett, C., Sun, M.T., Dolan, B.: Contextualized perturbation for textual adversarial attack. In: Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 5053–5069 (2021)
12. Li, J., Ji, S., Du, T., Li, B., Wang, T.: Textbugger: Generating adversarial text against real-world applications. In: 26th Annual Network and Distributed System Security Symposium. p. 55 (2019)

13. Li, L., Ma, R., Guo, Q., Xue, X., Qiu, X.: BERT-ATTACK: Adversarial attack against BERT using BERT. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing. pp. 6193–6202 (2020)
14. Ma, G., Shi, L., Guan, Z.: Adversarial text generation via probability determined word saliency. In: International Conference on Machine Learning for Cyber Security. pp. 562–571 (2020)
15. Maas, A., Daly, R.E., Pham, P.T., Huang, D., Ng, A.Y., Potts, C.: Learning word vectors for sentiment analysis. In: Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 142–150 (2011)
16. Miller, George, A.: Wordnet: a lexical database for english. In: Communications of the ACM. vol. 38, pp. 39–41 (1995)
17. Morris, J., Lifland, E., Yoo, J.Y., Grigsby, J., Jin, D., Qi, Y.: 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. pp. 119–126 (2020)
18. Mrksic, N., Séaghdha, D.Ó., Thomson, Blaise.and Gasic, M., Rojas-Barahona, Lina, M., Su, P.H., Vandyke, D., Wen, T.H., Young, S.J.: Counter-fitting word vectors to linguistic constraints. In: Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 142–148 (2016)
19. Pang, B., Lee, L.: Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In: Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 115–124 (2005)
20. Ren, S., Deng, Y., He, K., Che, W.: Generating natural language adversarial examples through probability weighted word saliency. In: Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 1085–1097 (2019)
21. Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C.D., Ng, A.Y., Potts, C.: Recursive deep models for semantic compositionality over a sentiment treebank. In: Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing. pp. 1631–1642 (2013)
22. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)
23. Wang, Z., Zheng, Y., Zhu, H., Yang, C., Chen, T.: Transferable adversarial examples can efficiently fool topic models. In: Computers & Security. vol. 118, p. 102749 (2022)
24. Williams, A., Nangia, N., Bowman, S.: A broad-coverage challenge corpus for sentence understanding through inference. In: Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 1112–1122 (2018)
25. Yang, X., Liu, W., Tao, D., Liu, W.: Besa: Bert-based simulated annealing for adversarial text attacks. In: International Joint Conference on Artificial Intelligence. pp. 3293–3299 (2021)
26. Zang, Y., Qi, F., Yang, C., Liu, Z., Zhang, M., Liu, Q., Sun, M.: Word-level textual adversarial attacking as combinatorial optimization. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 6066–6080 (2020)