

# BLOCK-SPARSE ADVERSARIAL ATTACK TO FOOL TRANSFORMER-BASED TEXT CLASSIFIERS

---

Sahar Sadrizadeh<sup>1</sup>, Ljiljana Dolamic<sup>2</sup>, and Pascal Frossard<sup>1</sup>

<sup>1</sup>École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

<sup>2</sup>Armasuisse S+T, Thun, Switzerland

Email: [sahar.sadrizadeh@epfl.ch](mailto:sahar.sadrizadeh@epfl.ch)

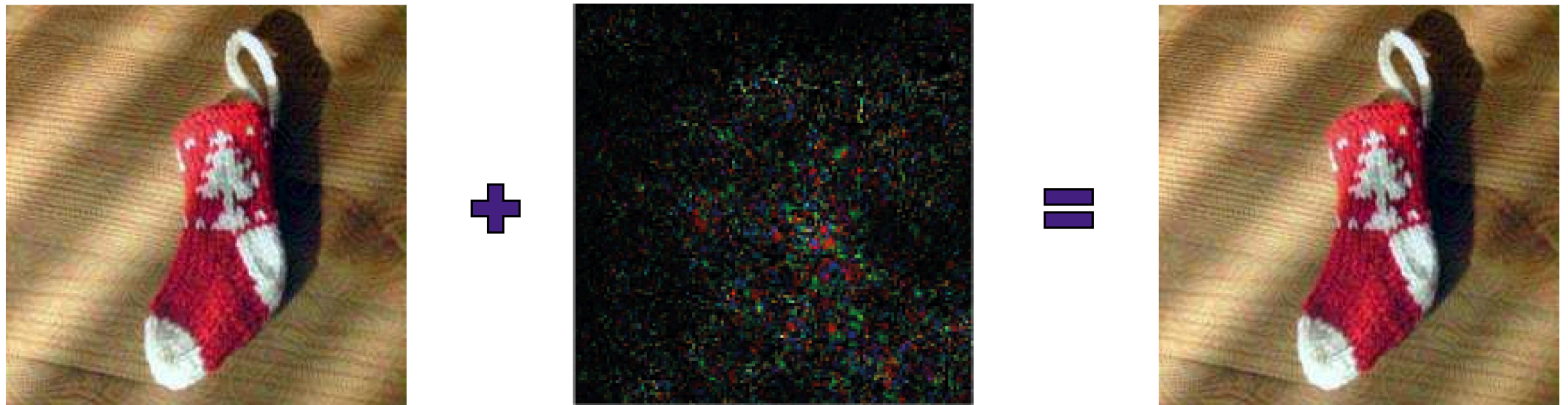


# Outline

- **Introduction**
- **Proposed Method**
- **Experiments**
- **Conclusion**

# Introduction

# Adversarial Attacks against Images



Socks

Imperceptible  
Perturbation

Indian Elephant

Moosavi-Dezfooli, et al. "DeepFool: a simple and accurate method to fool deep neural network." *CVPR*, 2016

# Challenges in NLP

- Input Data is **Discrete**:
  - The system is not end-to-end differentiable (optimizing over discrete input is difficult)
  - Working on embedding space may result in invalid text
- **Perturbation** of the adversarial example from the original sentence should be **imperceptible**:
  - Semantics: Similar meaning
  - Syntax: Correct grammar (fluent sentence)

Methods for adversarial attack against images is not applicable

# Word Substitution Adversarial Attacks

## 1. Rank Words

- White-box attack: Gradients
- Black-box attack: Masking each word and compare the output scores

## 2. Find possible **replacements** to maintain semantics

- Masked Language Models
- Embedding Space

## 3. **Substitute** with the word that changes the network result

- Two checks: **POS** and similar semantics (**Universal Sentence Encoder**)

Jin, et al. "Is bert really robust? a strong baseline for natural language attack on text classification and entailment." *AAAI*, 2020.

Li, et al. "BERT-ATTACK: Adversarial Attack Against BERT Using BERT." *EMNLP*, 2020.

# Problem Formulation

# White-box attack against transformers

Assumption: adversarial example and input have equal lengths

Original sentence:  $\mathbf{x} = x_1 x_2 \dots x_n \xrightarrow{\text{embedding}} \mathbf{e}_{\mathbf{x}} = [\text{emb}(x_1), \text{emb}(x_2), \dots, \text{emb}(x_n)]$   
 $y = f(\mathbf{x})$

Adversarial example:  $\mathbf{x}' = x'_1 x'_2 \dots x'_n \xrightarrow{\text{embedding}} \mathbf{e}_{\mathbf{x}'} = \mathbf{e}_{\mathbf{x}} + \mathbf{r}_{\mathbf{x}}$   
 $f(\mathbf{x}') \neq y$

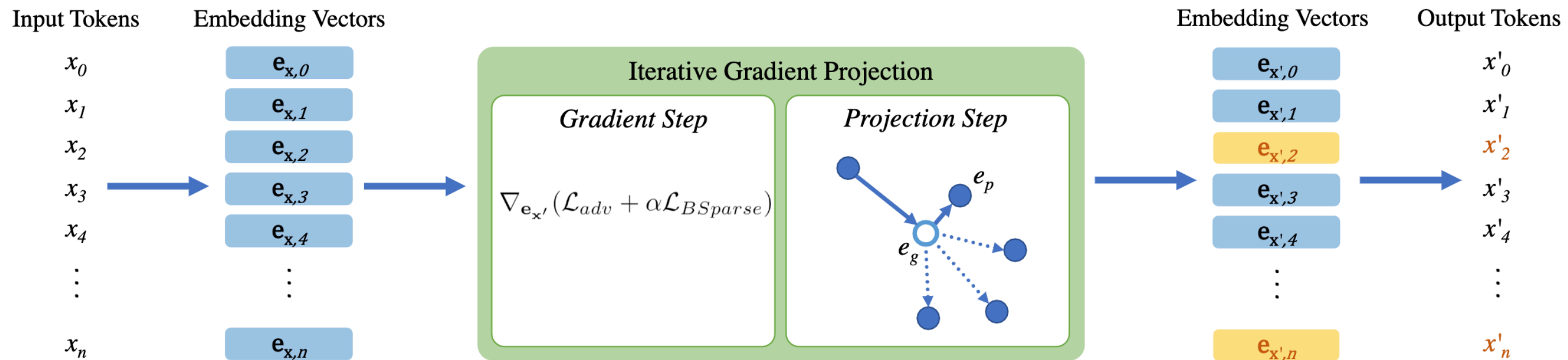
Perturbation:  $\mathbf{r}_{\mathbf{x}} = [\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_n] \xrightarrow{\text{Only few word changes}} \mathbf{r}_{\mathbf{x}} \text{ is block sparse}$

$$\mathcal{L}_{Adv} = -\mathcal{L}_f(\mathbf{e}_{\mathbf{x}'}, y) \qquad \mathcal{L}_{BSparse} = \sum_{i=1}^n \|\mathbf{r}_i\|_2$$

$$\hat{\mathbf{e}}_{\mathbf{x}'} = \underset{\mathbf{e}_{\mathbf{x}'} \in \mathcal{E}_{\mathcal{V}}}{\operatorname{argmin}} \mathcal{L}_{Adv} + \alpha \mathcal{L}_{BSparse}$$



# Block Diagram



Block diagram of the proposed method.

# Algorithm

Stop if the label is wrong

Gradient descent

Projection

---

## Algorithm 1 Block-Sparse Adversarial Attack

---

1: **Input:**

$f(\cdot)$ : Target classifier model,  $\mathcal{V}$ : Vocabulary set  
 $\mathbf{x}$ : Tokenized input sentence,  $lr$ : Learning rate  
 $A$ : Set of decreasing values for Hyper-parameter  $\alpha$  to control the importance of the block-sparsity term  
 $K$ : Maximum number of iterations

2: **Output:**

$\mathbf{x}'$ : Generated adversarial example

3: **procedure**

**initialization:**

4: **buffer**  $\leftarrow$  empty,  $y \leftarrow f(\mathbf{x})$ ,  $k \leftarrow 0$

5:  $\forall i \in \{1, \dots, n\} \quad \mathbf{e}_{g,i} \leftarrow \text{emb}(x_i)$

6: **for**  $\alpha$  in  $A$  **do**

7: **while**  $f(\mathbf{e}_p) = y$  and  $k \leq K$  **do**

8:  $k \leftarrow k + 1$

**Step 1:** Gradient descent in the continuous embedding space:

9:  $\mathbf{e}_g \leftarrow \mathbf{e}_g - lr \cdot \nabla_{\mathbf{e}_{x'}} (\mathcal{L}_{adv} + \alpha \mathcal{L}_{BSparse})$

**Step 2:** Projection to the discrete subspace  $\mathcal{E}_{\mathcal{V}}$  and update if the sentence is new:

10: **for**  $i \in \{1, \dots, n\}$  **do**

11:  $\mathbf{e}_{p,i} \leftarrow \underset{\mathbf{e} \in \mathcal{E}_{\mathcal{V}}}{\text{argmin}} \frac{\mathbf{e}^\top \mathbf{e}_{g,i}}{\|\mathbf{e}\|_2 \cdot \|\mathbf{e}_{g,i}\|_2}$

12: **end for**

13: **if**  $\mathbf{e}_p$  not in **buffer** **then**

14: **add**  $\mathbf{e}_p$  to **buffer**

15:  $\mathbf{e}_g \leftarrow \mathbf{e}_p$

16: **end if**

17: **end while**

18: **if**  $f(\mathbf{e}_p) \neq y$  **then**

19: **break** (adversarial example is found)

20: **end if**

21: **end for**

22: **return**  $\mathbf{e}_{x'} \leftarrow \mathbf{e}_p$

23: **end procedure**

---

# Experiments

# Experimental Setup

- ▶ **Target model:** GPT-2 Transformer
- ▶ **3 Datasets:**
  - MNLI (natural language inference),
  - AG News (news categorization)
  - Yelp Reviews (sentiment classification).
- ▶ **Baseline:** GBDA (optimization-based white-box attack against transformers)

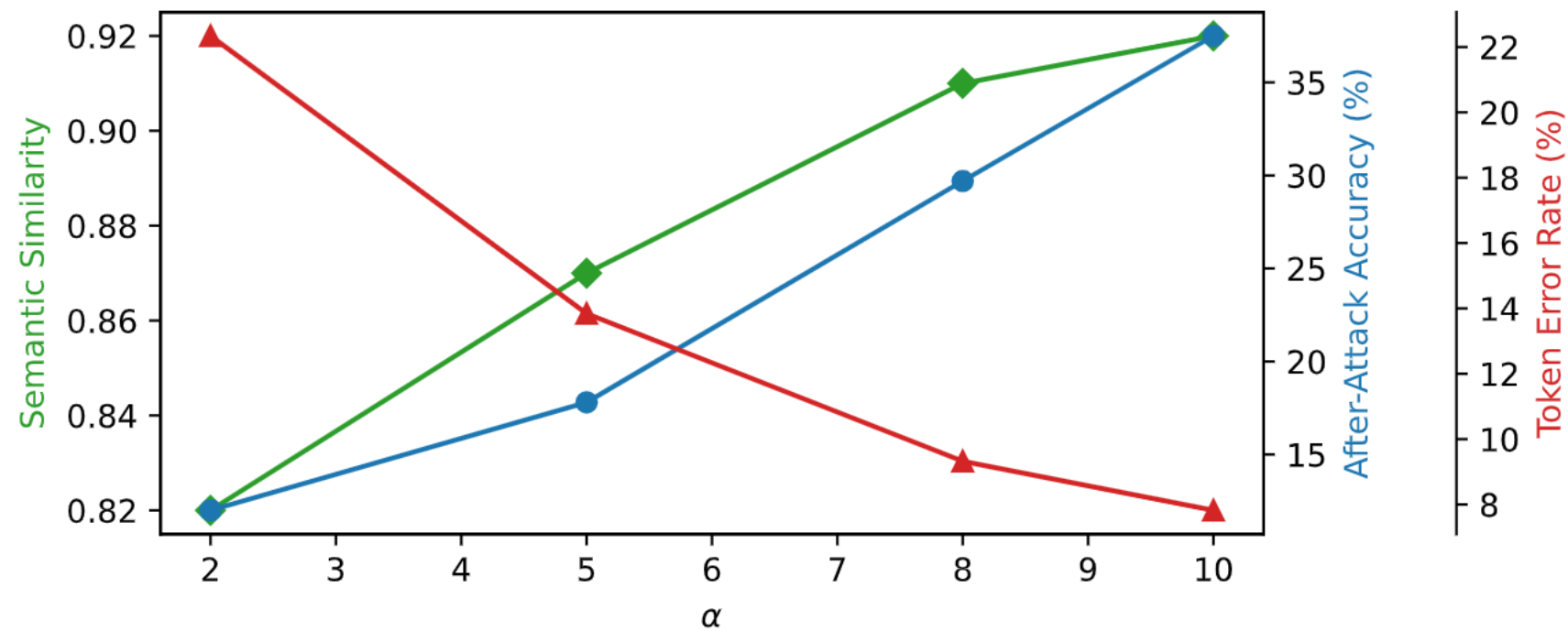
Guo, et al. "Gradient-based Adversarial Attacks against Text Transformers." *EMNLP*, 2021.

# Experimental Results

Performance of white-box attack against fine-tuned GPT-2

Method	Ag News		MNLI		Yelp Reviews	
	Adv. Acc.	Sim.	Adv. Acc.	Sim.	Adv. Acc.	Sim.
Proposed	0.4	0.87	3.0 (1.3)	0.85 (0.82)	1.8	0.87
GBDA	6.6	0.90	2.8 (11.0)	0.82 (0.88)	2.9	0.94

# Ablation Study



Effect of the hyper-parameter  $\alpha$  on the performance of our attack.



# Examples

Dataset	Sentence	Prediction	Text
MNLI	Original	Neutral (97.26%)	Premise: In the summer, the Sultan's Pool, a vast outdoor amphitheatre, stages rock concerts or other big-name events. Hypothesis: <b>Most</b> rock concerts take place in the Sultan's Pool amphitheatre.
	Adversarial	<b>Entailment (99.19%)</b>	Premise: In the summer, the Sultan's Pool, a vast outdoor amphitheatre, stages rock concerts or other big-name events. Hypothesis: <b>Many</b> rock concerts take place in the Sultan's Pool amphitheatre.
AG News	Original	Sci/Tech (99.39%)	Motorola and HP in <b>Linux</b> tie-up Motorola plans to sell mobile phone network equipment that uses <b>Linux</b> -based code, a step forward in network gear makers #39; efforts to rally around a standard.
	Adversarial	<b>Business (83.56%)</b>	Motorola and HP in <b>PC</b> tie-up Motorola plans to sell mobile phone network equipment that uses <b>PC</b> -based code. a step forward in network gear makers #39; efforts to rally around a standard.
Yelp	Original	Negative (99.90%)	This place holds a nostalgic appeal for people born and raised in Pittsburgh who grew up eating here. If that experience is what your looking for, please visit. If you're looking for a tasty meal, go somewhere else. 5 stars for history, <b>0</b> for food quality and flavor.
	Adversarial	<b>Positive (96.54%)</b>	This place holds a nostalgic appeal for people born and raised in Pittsburgh who grew up eating here. If that experience is what your looking for, please visit. If you're looking for a tasty meal, go somewhere else. 5 stars for history, <b>1</b> for food quality and flavor.

# Conclusion

- ▶ Optimization problem to generate adversarial examples
  - Block sparsity constraint to ensure few tokens are modified
  - Solve by Gradient projection
- ▶ *Comparable* Performance with GBDA
  - Drops the classification accuracy to *less than 5%*
  - Semantic similarity is *more than 80%*
- ▶ **Source code:** <https://github.com/ssadrizadeh/transformer-text-classifier-attack>



# Thanks for Your Attention



[Source Code]

Sahar Sadrizadeh

`sahar.sadrizadeh@epfl.ch`

Signal Processing Laboratory (LTS4), EPFL

**EPFL**

 **icassp 2022**  
*Singapore*