# Adversarial Examples Against a BERT ABSA Model

## Fooling Bert With L33T, Misspellign, and Punctuation,
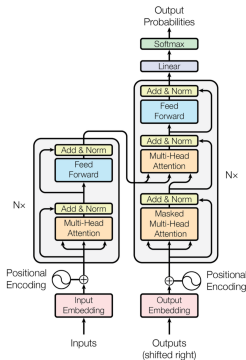
**N. Hofer, P. Schöttle, A. Rietzler, S. Stabinger**
**August, 2021**

Bidirectional Encoder Representations from Transformers - BERT



Transformer Model Architecture (Vaswani et al., 2017)

$$x$$
"panda"
57.7% confidence

$$+ .007 \times$$

$$\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"nematode"
8.2% confidence

$$=$$

$$\boldsymbol{x} + \epsilon\,\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"gibbon"
99.3 % confidence

Adversarial Examples in Computer Vision (Goodfellow et al, ICLR 2015)

Object detection in autonomous driving (Source: becominghuman.ai)

# Motivation



Tweet containing misleading information regarding Covid-19, detected and labeled correctly

# Motivation



Tweet containing misleading information regarding Covid-19, detected and labeled correctly

# Motivation



Tweet containing misleading information regarding Covid-19. Potential problems due to the use of Leet Speak.

# Adversarial Attacks

*"We're starting to deploy machine learning as a technology that can fail so we need to have some checks in place."* Nicolas Papernot, 2017

# Adversarial Attacks against BERT for ABSA

1. Fine-Tuning BERT base for ABSA

# Adversarial Attacks against BERT for ABSA

1. Fine-Tuning BERT base for ABSA

# Adversarial Attacks against BERT for ABSA

| 1. Fine-Tuning BERT base for ABSA |
|---|

Aspect-based Sentiment Analysis

# Adversarial Attacks against BERT for ABSA

1. Fine-Tuning BERT base for ABSA

Aspect-based Sentiment Analysis

*The computer is excellent for gaming but I think it is way too expensive!!*

# Adversarial Attacks against BERT for ABSA

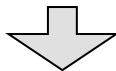| 1. Fine-Tuning BERT base for ABSA |
|---|

Aspect-based Sentiment Analysis

*The computer is* `excellent for gaming` *but I think it is* `way too expensive` *!!*

Aspect: Gaming, Sentiment: POS
Aspect: Price, Sentiment: NEG

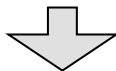# Adversarial Attacks against BERT for ABSA

1. Fine-Tuning BERT base for ABSA

2. Identify Important Word

# Adversarial Attacks against BERT for ABSA
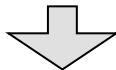
1. Fine-Tuning BERT base for ABSA

2. Identify Important Word

Leave-One-Out Method

# Adversarial Attacks against BERT for ABSA

| 1. Fine-Tuning BERT base for ABSA |
| :---: |

$$\Downarrow$$

| 2. Identify Important Word |
| :---: |

## Leave-One-Out Method

**The computer is** `excellent for gaming` **but I think it is** `way too expensive!!` Gaming - POS; Price - NEG

**The** - *computer is* `excellent for gaming` *but I think it is* `way too expensive!!` Gaming - POS; Price - NEG

**computer** - *The is* `excellent for gaming` *but I think it is* `way too expensive!!` Gaming - POS; Price - NEG
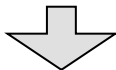
**is** - *The computer* `excellent for gaming` *but I think it is* `way too expensive!!` Gaming - POS; Price - NEG

**excellent** - *The computer is for gaming but I think it is* `way too expensive!!` Price - NEG
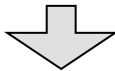
**for** - *The computer is* `excellent gaming` *but I think it is* `way too expensive!!` Gaming - POS; Price - NEG

# Adversarial Attacks against BERT for ABSA

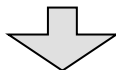1. Fine-Tuning BERT base for ABSA

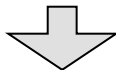2. Identify Important Word

3. Modification of Important Words

# Adversarial Attacks against BERT for ABSA

1. Fine-Tuning BERT base for ABSA

2. Identify Important Word

3. Modification of Important Words

| Leetspeak | Misspellings | Punctuation |

# Adversarial Attacks

## Objectives

- Semantic Meaning
- Inconspicousness
- Relevance

# Adversarial Attacks against BERT for ABSA

## 1. Leetspeak

*The computer is* `excellent for gaming` *but I think it is* `way too expensive!!`

Aspect: Gaming, Sentiment: POS
Aspect: Price, Sentiment: NEG

*Original important word: **excellent***
*Modified important word:* ***exce11ent***

*The computer is* `exce11ent for gaming` *but I think it is* `way too expensive!!`

Aspect: Gaming, Sentiment: NEG
Aspect: Price, Sentiment: NEG

# Adversarial Attacks against BERT for ABSA

| 2. Misspellings |
|---|

*The computer is* **excellent for gaming** *but I think it is* **way too expensive!!**

Aspect: Gaming, Sentiment: POS
Aspect: Price, Sentiment: NEG

*Original important word:* **excellent**
*Modified important word:* **ecxellent**

*The computer is ecxellent for gaming but I think it is* **way too expensive!!**

Aspect: Price, Sentiment: NEG

# Adversarial Attacks against BERT for ABSA

## 3. Punctuation

*The computer is* `excellent for gaming` *but I think it is* `way too expensive!!`

Aspect: Gaming, Sentiment: POS
Aspect: Price, Sentiment: NEG

*Original important word: **excellent***
*Modified important word: **excellent,***

*The* `computer is excellent,` `for gaming but I think` *it is* `way too expensive!!`

Aspect: Laptop (general), Sentiment: NEG
Aspect: Gaming, Sentiment: NEG
Aspect: Price, Sentiment: NEG

tba

# Conclusion

*Our reliance on deep learning based language models for real-world (security-relevant) applications is questionable.*

# Adversarial Examples Against a BERT ABSA Model

## Fooling Bert With L33T, Misspellign, and Punctuation,

**N. Hofer, P. Schöttle, A. Rietzler, S. Stabinger**
**August, 2021**