Paper Seminar

# FreeLB: Enhanced Adversarial Training for Natural Language Understanding

**Zhu et al., 2020, ICLR**

**Myeongsup Kim**

Integrated M.S./Ph.D. Student
Data Science & Business Analytics Lab.
School of Industrial Management Engineering
Korea University

Myeongsup_kim@korea.ac.kr

# Introduction

- Transformer-Based Language Model

# \<What This Seminar Does Not Cover\>

- ## Details of Transformer

  Vaswani et al., Attention is All You Need, NIPS, 2017

- ## Details of BERT and RoBERTa

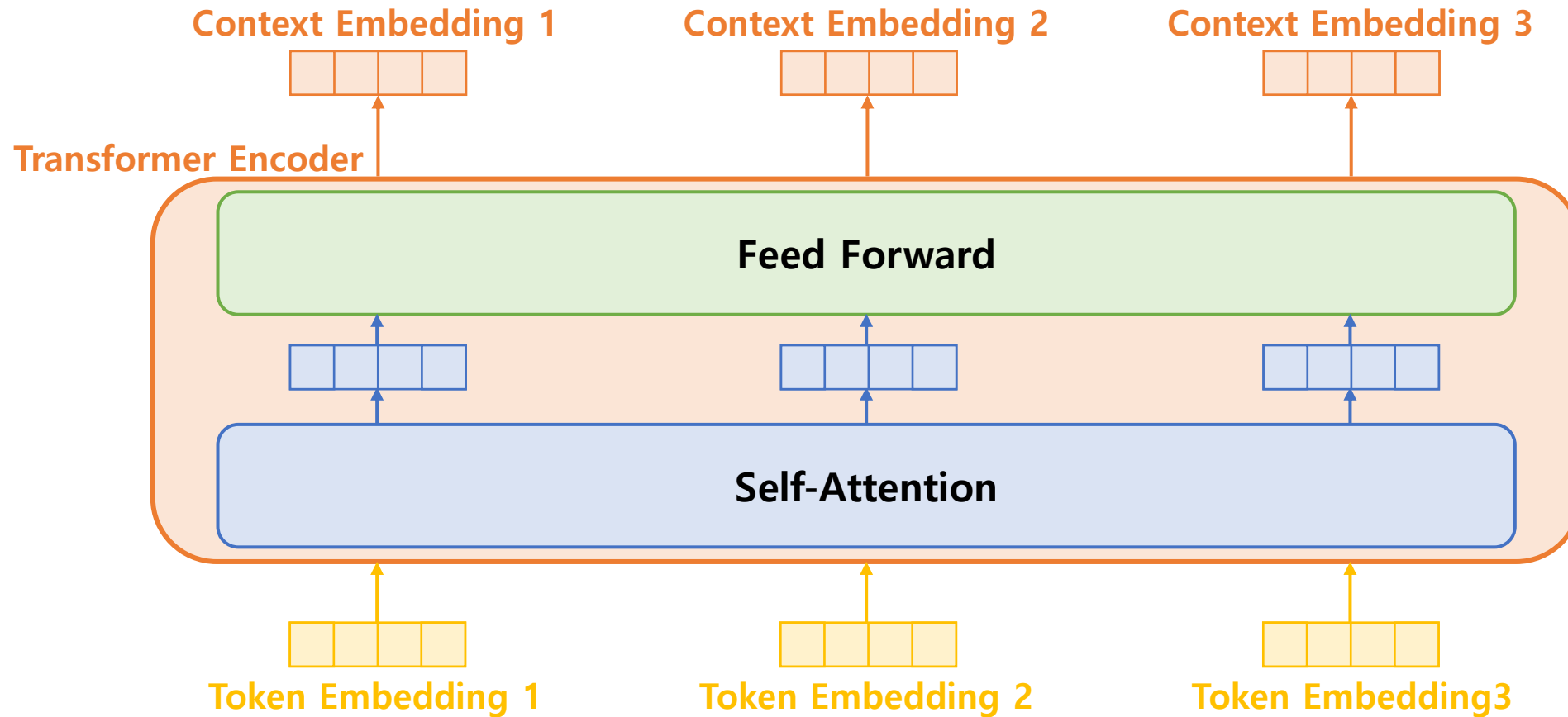  Devlin et al., BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, NAACL, 2019

  Liu et al., RoBERTa: A Robustly Optimized BERT Pretraining Approach, arXiv, 2019
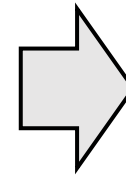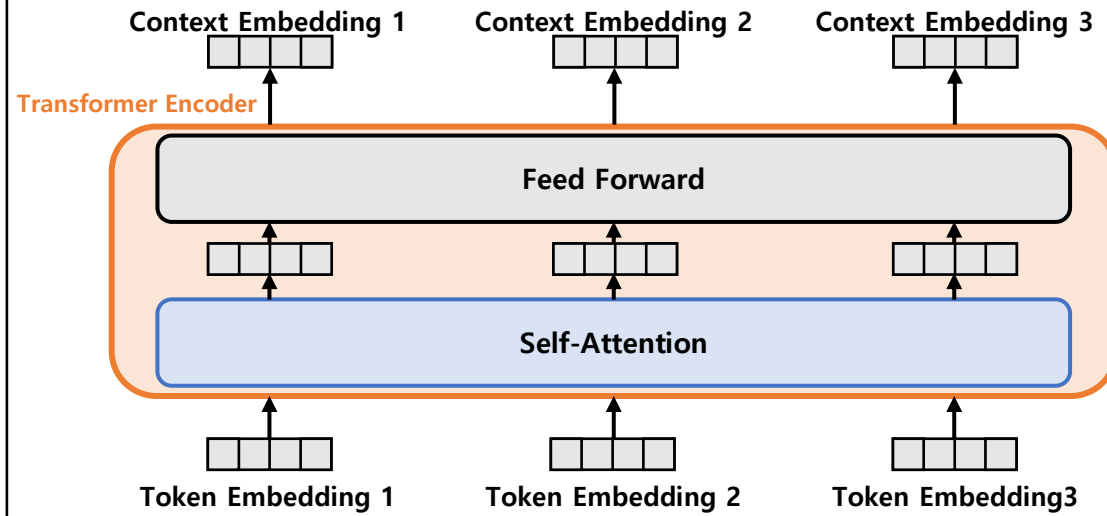
# Introduction
-Transformer-Based Language Model

# \<Transformer Encoder\>

**Context Embedding 1**

**Context Embedding 2**

**Context Embedding 3**

**Transformer Encoder**

**Feed Forward**

**Self-Attention**

**Token Embedding 1**

**Token Embedding 2**

**Token Embedding3**

# Introduction
## -Transformer-Based Language Model

# \<Self-Attention\>

# Introduction
-Transformer-Based Language Model

# <Self-Attention>



**Transformer Encoder**

Context Embedding 1    Context Embedding 2    Context Embedding 3

Feed Forward

Self-Attention

Token Embedding 1    Token Embedding 2    Token Embedding3

$$\frac{Q \times K^T}{\sqrt{d_k}} =$$

Q    K$^T$    Attention Score Matrix

$$softmax(\quad) \times \quad =$$

Attention Score Matrix    V    Attention Value Matrix

$$concat(\quad) \times \quad =$$

Attention Value Matrix    W    Context Vector

# Introduction
**-Transformer-Based Language Model**

## <Self-Attention>

# Introduction
-Transformer-Based Language Model

# <Self-Attention>

# <Self-Attention>

# Introduction
-Transformer-Based Language Model

## <Self-Attention>

# Introduction
**-Transformer-Based Language Model**

# <Self-Attention>

Context Embedding 1     Context Embedding 2     Context Embedding 3

**Transformer Encoder**

**Feed Forward**

**Self-Attention**

Token Embedding 1     Token Embedding 2     Token Embedding3

Self-Attended Token Embedding 1

Self-Attended Token Embedding 2

Self-Attended Token Embedding 3

# Introduction
-Transformer-Based Language Model

## <Feed Forward>

# Introduction
**-Transformer-Based Language Model**

## <Feed Forward>



Transformer Encoder

Context Embedding 1    Context Embedding 2    Context Embedding 3

Feed Forward

Self-Attention

Token Embedding 1    Token Embedding 2    Token Embedding3

Context Embedding 3    Context Embedding 1

Context Embedding 2

# <Contextualized Representation>

# <Contextualized Representation>



German article "die"

Was der Fall ist, **die** Tatsache, ist das Bestehen von Sachverhalten.

über **die** Verhandlungen der Königl.

single person dies ⟷ multiple people die

Chernenko became the first Soviet leader to **die** in less than three years

Vaughan's ultimate fantasy was to **die** in a head-on collision with movie star Elizabeth Taylor

Over 60 people **die** and over 100 are unaccounted for.

Many more **die** from radiation sickness, starvation and cold.

a playing die

Players must always move a token according to the **die** value

The faces of a **die** may be placed clockwise or counterclockwise

**<Embeddings for the Word "die" in Different Contexts>**

# <Pre-Training>

Pre-Training Task

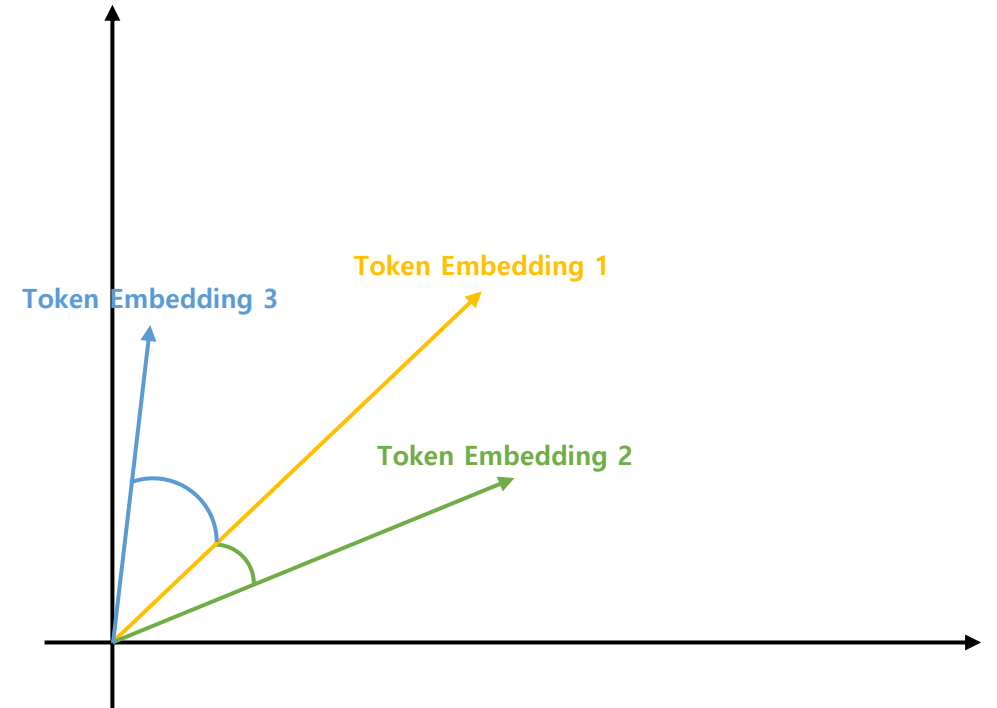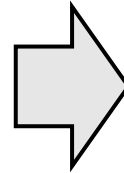Context Embedding 1    Context Embedding 2    Context Embedding 3
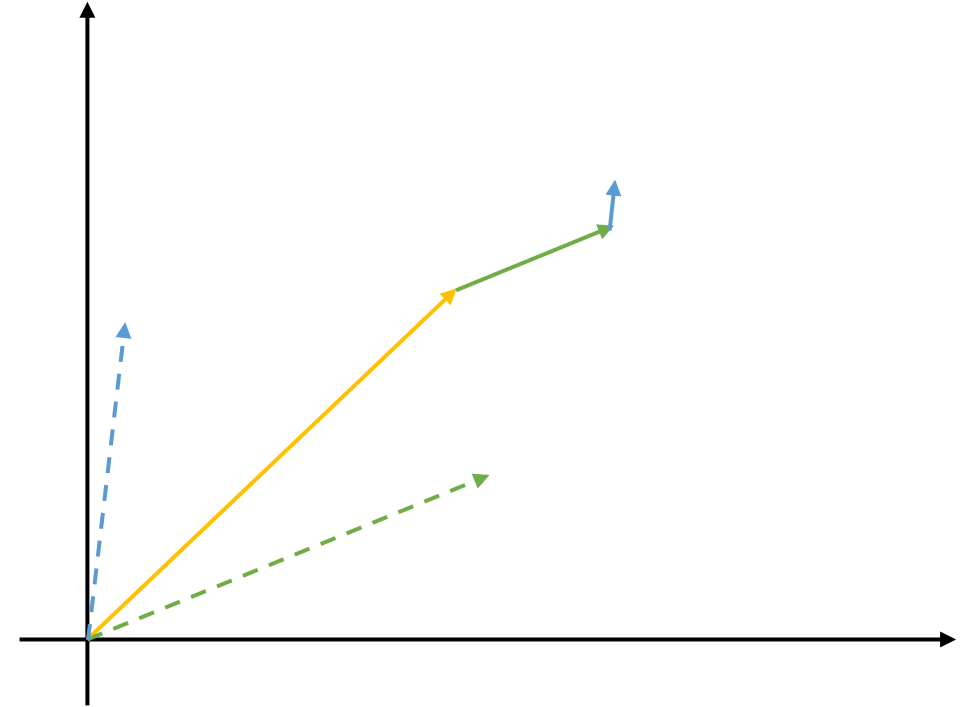
Encoder

Very Large Text Corpora

# Introduction
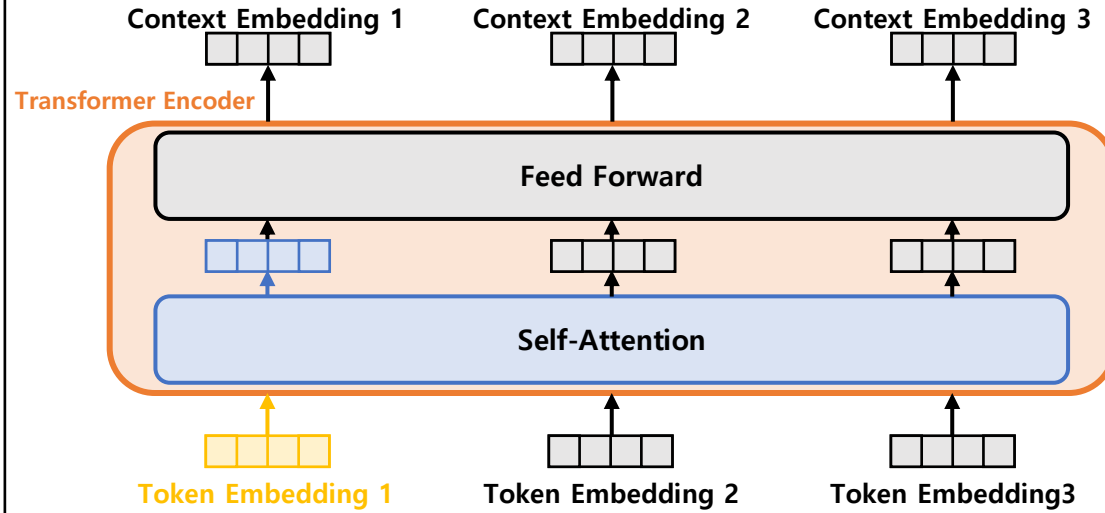-Transformer-Based Language Model

## <Fine-Tuning>

Ground Truth of
Specific Task

Simple Classifier (Linear or Small MLP)

Context Embedding 1    Context Embedding 2    Context Embedding 3

## Encoder

Small Task Specific Data

# <Two Branches of Language Model Research>

**"Bigger, Larger, Stronger"**　　　**"Small, But Better Performance"**

# <Two Branches of Language Model Research>

## "Bigger, Larger, Stronger"

- Training Deep and Large Models with Huge Data

## "Small, But Better Performance"

- Improving Performance without Changing the Structure of the Language Model

- Changing the Structure of the Model Without Significantly Increasing the Parameters

# <Two Branches of Language Model Research>

## "Bigger, Larger, Stronger"

- **Text To Text Transfer Transformer (T5)**
  - ✓ 11B Parameters
  - ✓ State-of-the-art in **GLUE**, etc.

  Raffel et al., Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer, JMLR, 2020

- **Generative Pre-Trained Transformer 3 (GPT-3)**
  - ✓ 175B Parameters
  - ✓ State-of-the-art in Many Benchmarks with Zero/Few Shot Setting

  Brown et al., Language Models are Few-Shot Learners, NeurIPS, 2020

## "Small, But Better Performance"

- **SMART**
  - ✓ 356M Parameters
  - ✓ Beat T5 in 3 Tasks of GLUE

  Jiang et al., SMART: Robust and Efficient Fine-Tuning for Pre-trained Natural Language Models through Principled Regularized Optimization, ACL, 2020

- **Pattern-Exploiting Training (PET)**
  - ✓ 223M Parameters
  - ✓ Beat GPT-3 in SuperGLUE with Few Shot Setting

  Schick and Schutze, It's Not Just Size That Matters: Small Language Models are Also Few-Shot Learners, arXiv, 2020

# <Two Branches of Language Model Research>

| **"Bigger, Larger, Stronger"** | **"Small, But Better Performance"** |
|---|---|

- Training Deep and Large Models with Huge Data

- Improving Performance **without Changing the Structure** of the Language Model

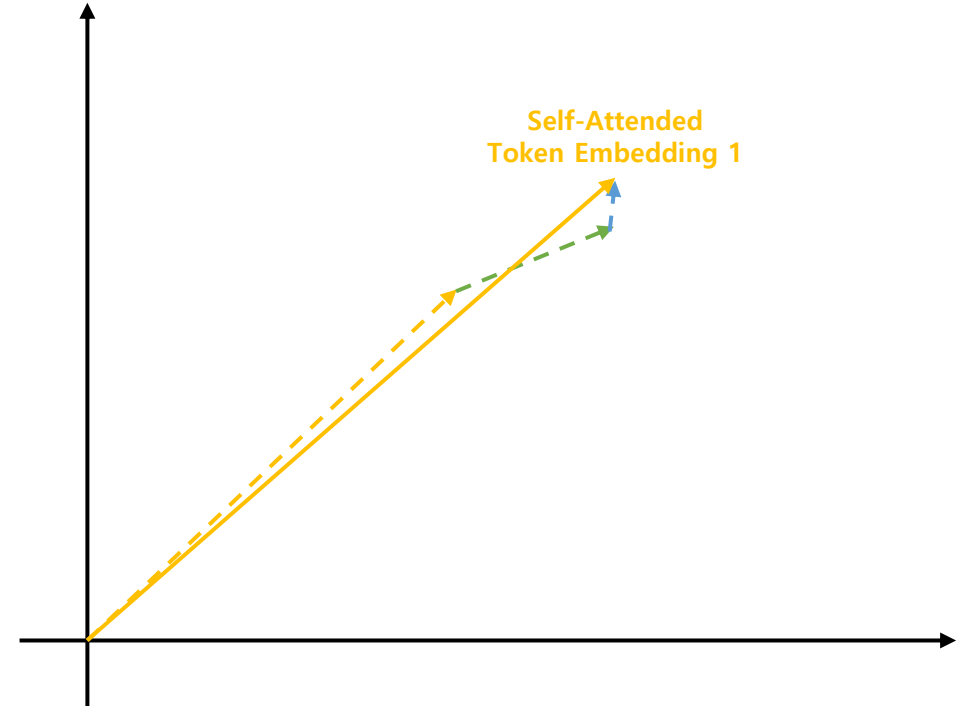- Changing the Structure of the Model Without Significantly Increasing the Parameters
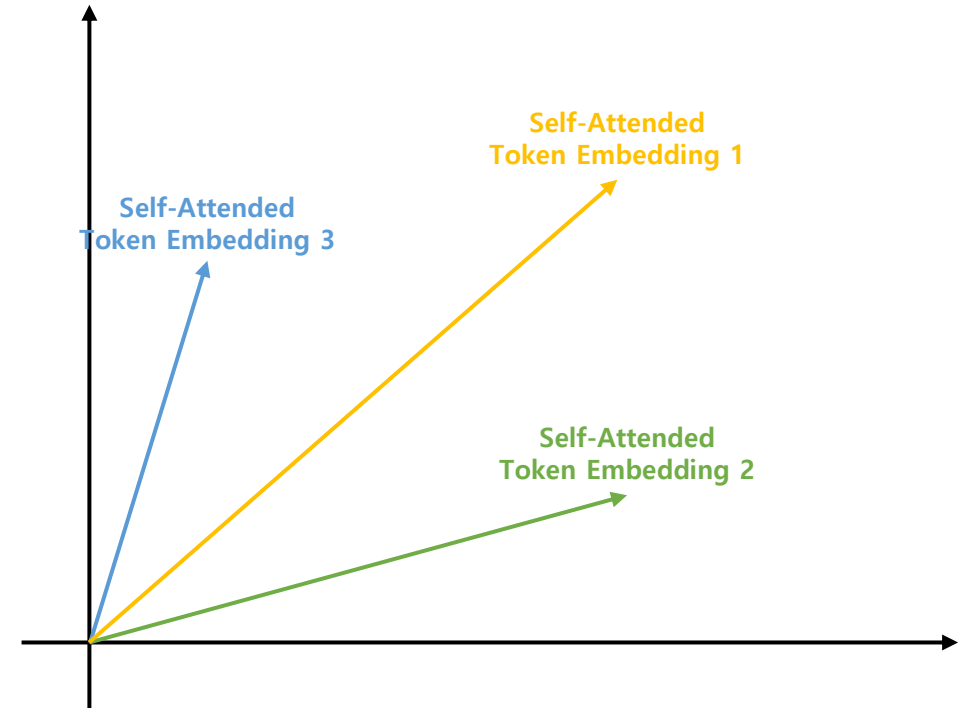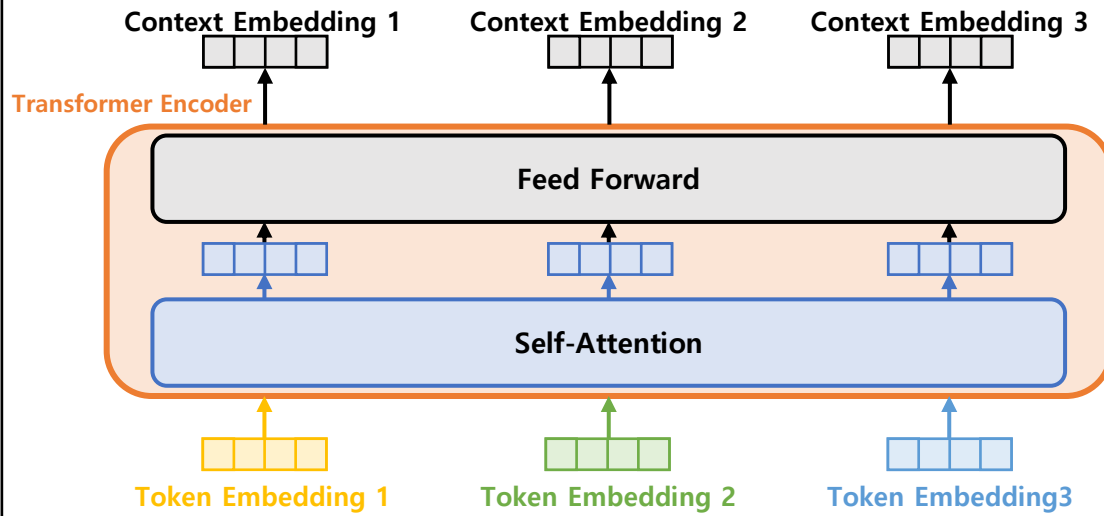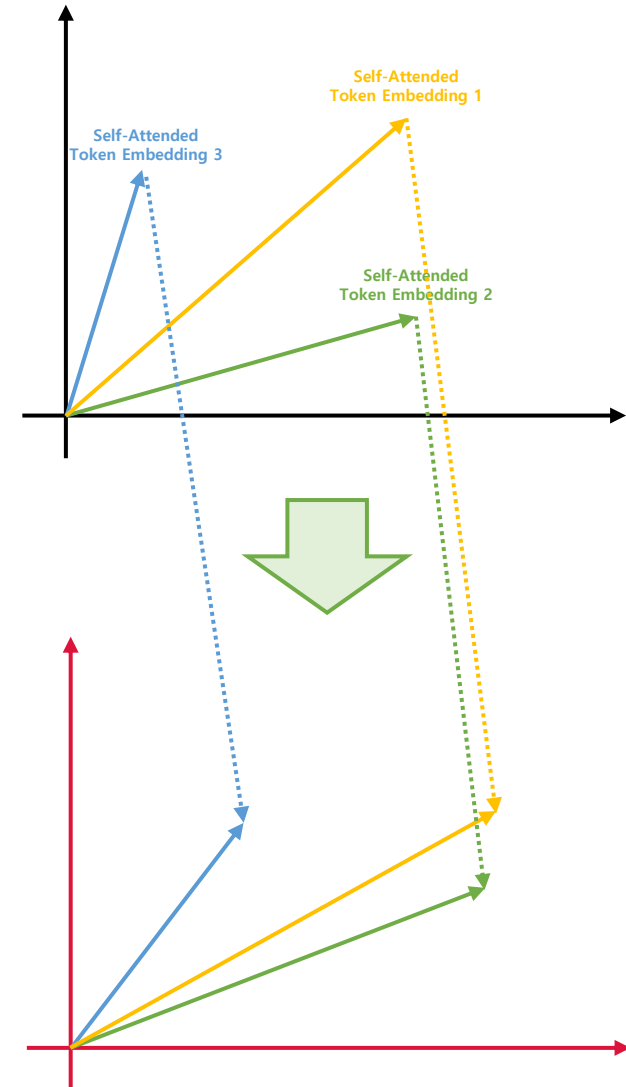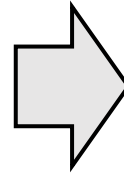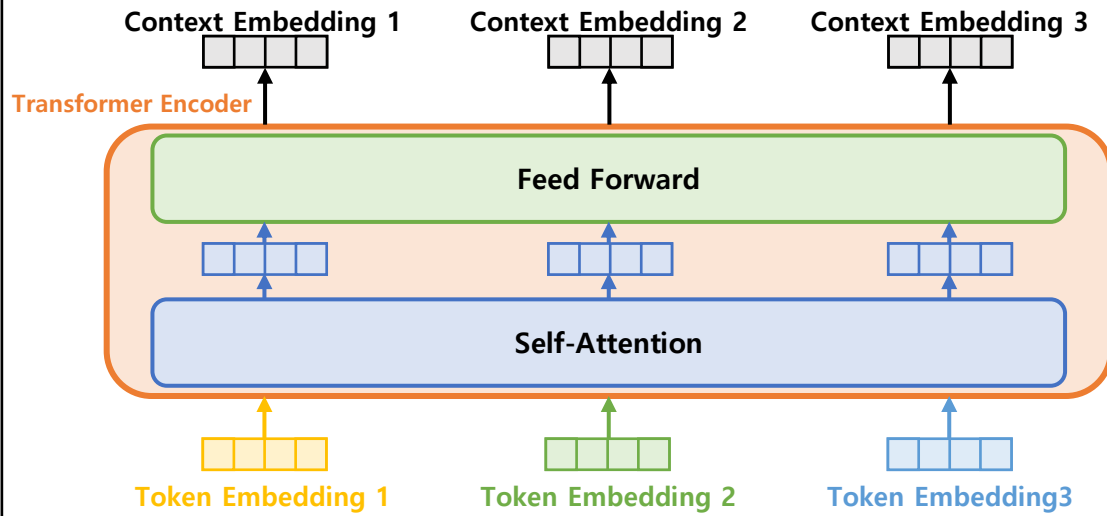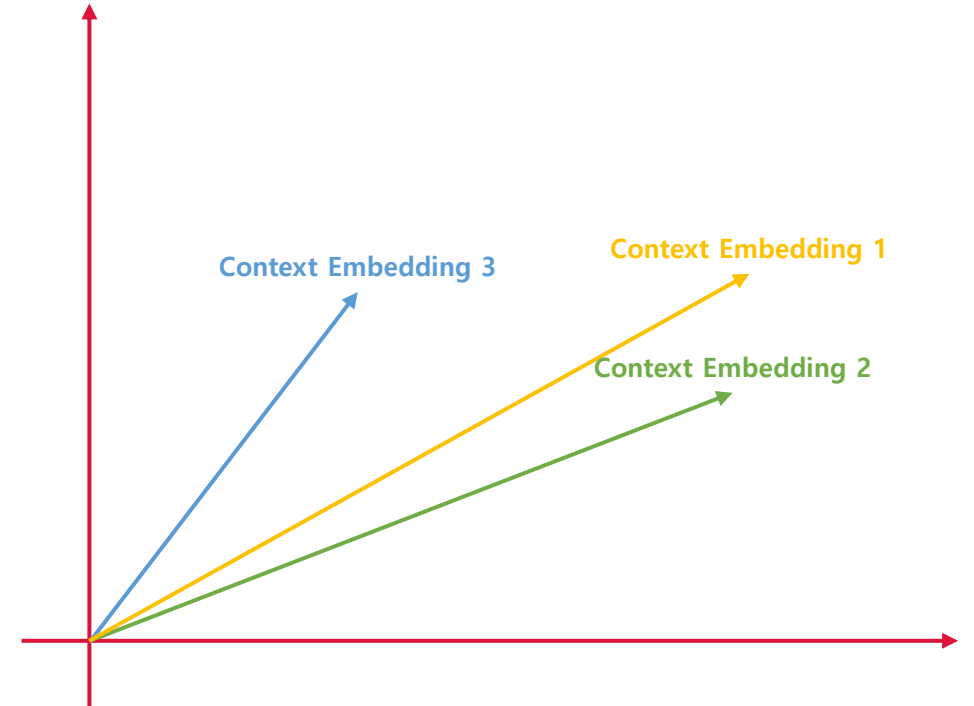
# Introduction
**-Transformer-Based Language Model**

## \<Task-Adaptive Pre-Training\>

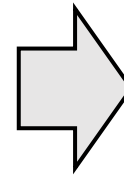| Small Task Specific Data | ⊄ | Very Large Text Corpora |
|---|---|---|

# Introduction
**-Transformer-Based Language Model**

&lt;Task-Adaptive Pre-Training&gt;

# Introduction

**-Transformer-Based Language Model**

## <Text Augmentation>

Small Task Specific Data → Augmented Task Specific Data
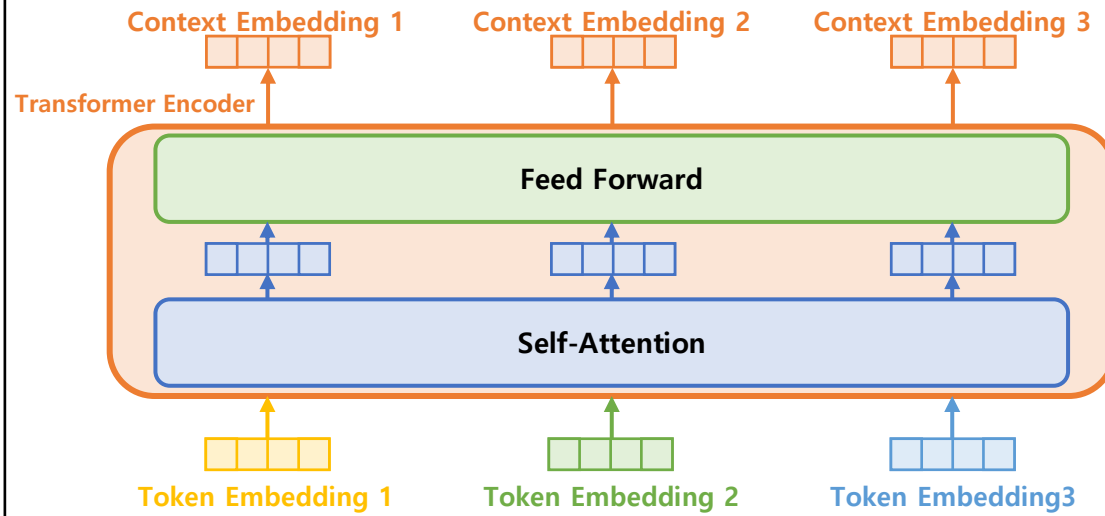
# Introduction
## -Transformer-Based Language Model

## \<Back Translation\>

<Text Augmentation>

Ground Truth of Specific Task

Simple Classifier (Linear or Small MLP)

Context Embedding 1    Context Embedding 2    Context Embedding 3

Encoder

Small Task Specific Data → Augmented Task Specific Data

# Introduction
-Transformer-Based Language Model

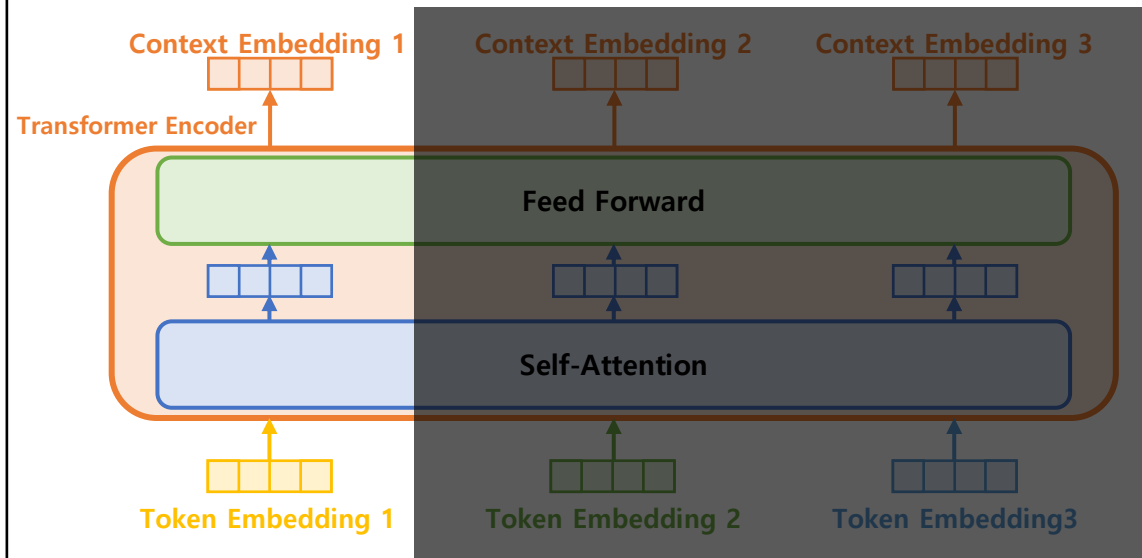## <Generalization Performance>

Context Embedding 1     Context Embedding 2     Context Embedding 3

Transformer Encoder

**Encoder**

Token Embedding 1     Token Embedding 2     Token Embedding3

Token 1     Token 2     Token 3

**Generalization Performance in Fine-Tuning**

Context Embedding

**Task Adaptation**

**Text Augmentation**

# &lt;Generalization Performance&gt;

# Introduction
**-Transformer-Based Language Model**

## \<Generalization Performance\>



**Context Embedding 1**

**Context Embedding 2**

**Context Embedding 3**

Transformer Encoder

**Encoder**

**Token Embedding 1**

**Token Embedding 2**

**Token Embedding3**

**Token 4**

**Token 2**

**Token 3**

Task Adaptation

**Text Augmentation**

**Generalization Performance in Fine-Tuning**

Context Embedding

# <Generalization Performance>



Context Embedding 1     Context Embedding 2     Context Embedding 3

Transformer Encoder

**Encoder**

Token Embedding 1     Token Embedding 2     Token Embedding3

Token 1     Token 2     Token 3

Generalization Performance
in Fine-Tuning

Context Embedding

**Embedding
Manipulation**

# Introduction
## -Transformer-Based Language Model

### &lt;Generalization Performance&gt;

Context Embedding 1    Context Embedding 2    Context Embedding 3

Transformer Encoder

**Encoder**

Token Embedding 1    Token Embedding 2    Token Embedding3

Token 1    Token 2    Token 3

Generalization Performance
in Fine-Tuning

Context Embedding

Embedding
Manipulation

**Can We Improve the Generalization Performance of the Language Model
by Manipulating Embeddings rather than Tokens?**

# FreeLB: Enhanced Adversarial Training for Natural Language Understanding

*Zhu et al., 2020, ICLR*

# FreeLB: Enhanced Adversarial Training for Natural Language Understanding

*Zhu et al., 2020, ICLR*

# Pre-requisites

- Adversarial Training

# Pre-requisites
**- Adversarial Training**

# \<Recap: Previous Seminar\>



Adversarial example

- 특수한 noise를 원본 example에 더하여 사람이 판단하기에는 똑같지만, machine이 판단하기에는 다른 class가 되는 example

- 예를 들어 오른쪽 이미지는 우리가 보기에는 여전히 고양이지만 DNN이 보기에는 오븐이 된다!

# Pre-requisites
**- Adversarial Training**

## <Adversarial Example>

$$+ .007 \times$$

$$=$$

$$x$$
"panda"
57.7% confidence

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

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

# Pre-requisites
**- Adversarial Training**

# \<Adversarial Training>



**\<Input Image>**

**Model**

**\<Training>**

**[Cat]**

**[Dog]**

**\<Classification>**

# <Adversarial Training>



<Perturbation>          <Input Image>

<Adversarial Example>

Model

<Inference>

[Cat]

[Dog]

<Classification>

# \<Adversarial Training\>



\<Adversarial Example\>

\<Inference\>

\<Classification\>

[Cat]

[Dog]

# Pre-requisites
**- Adversarial Training**

## \<Adversarial Training\>



\<Input Image\>

\<Perturbation Image\>

\<Adversarial Example\>

**Model**

\<Training\>

[Cat]

[Dog]

\<Classification\>

# Pre-requisites

**- Adversarial Training**

## <Adversarial Training>

# Pre-requisites
**- Adversarial Training**

## <Adversarial Training>

# Pre-requisites
**- Adversarial Training**

## \<Adversarial Training\>



**Cat**

**Dog**

# Pre-requisites
**- Adversarial Training**

## <Adversarial Training>



Cat

Dog

# Pre-requisites

**- Adversarial Training**

## \<Adversarial Training\>

# Pre-requisites
**- Adversarial Training**

## <Adversarial Training>



Cat

Dog

# Pre-requisites
**-Adversarial Training**

## \<Adversarial Training\>



?

# Adversarial Training for NLU

- PGD-Based Adversarial Training

- Large-Batch Adversarial Training for Free

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z,\,y)\sim\mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_\theta(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) \,/\, ||g(\delta_t)||_F)$$

$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_\theta(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$

Simple Classifier (Linear or Small MLP)

Context Embedding 1    Context Embedding 2    Context Embedding 3

Encoder

Embedding 1    Embedding 2    Embedding 3

Embedding

One-Hot Encoding of Input Data

**49/98**

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \left|\left|g(\delta_t)\right|\right|_F)$$

$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_{\theta}(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$

**Simple Classifier (Linear or Small MLP)**

Context Embedding 1   Context Embedding 2   Context Embedding 3

**Encoder**

Embedding 1   Embedding 2   Embedding 3

**Embedding**

**One-Hot Encoding of Input Data**

# \<Projected Gradient Descent\>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \left|\left|g(\delta_t)\right|\right|_F)$$

$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_{\theta}(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$

Simple Classifier (Linear or Small MLP)

Context Embedding 1    Context Embedding 2    Context Embedding 3

Encoder

Embedding 1    Embedding 2    Embedding 3

Embedding

One-Hot Encoding of Input Data

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(\boldsymbol{X} + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon} (\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$

$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_{\theta}(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$

# Adversarial Training for NLU
- **PGD-Based Adversarial Training**

## <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon} (\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$

$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_{\theta}(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$

**Simple Classifier (Linear or Small MLP)**

Context Embedding 1    Context Embedding 2    Context Embedding 3

**Encoder**

Embedding 1    Embedding 2    Embedding 3

**Embedding**

**One-Hot Encoding of Input Data**

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \left|\left| g(\delta_t) \right|\right|_F)$$

$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_{\theta}(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$

**Simple Classifier (Linear or Small MLP)**

Context Embedding 1    Context Embedding 2    Context Embedding 3

**Encoder**

Embedding 1    Embedding 2    Embedding 3

**Embedding**

**One-Hot Encoding of Input Data**

# \<Projected Gradient Descent\>

$$\min_{\theta} \mathbb{E}_{(Z, \, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon} (\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$
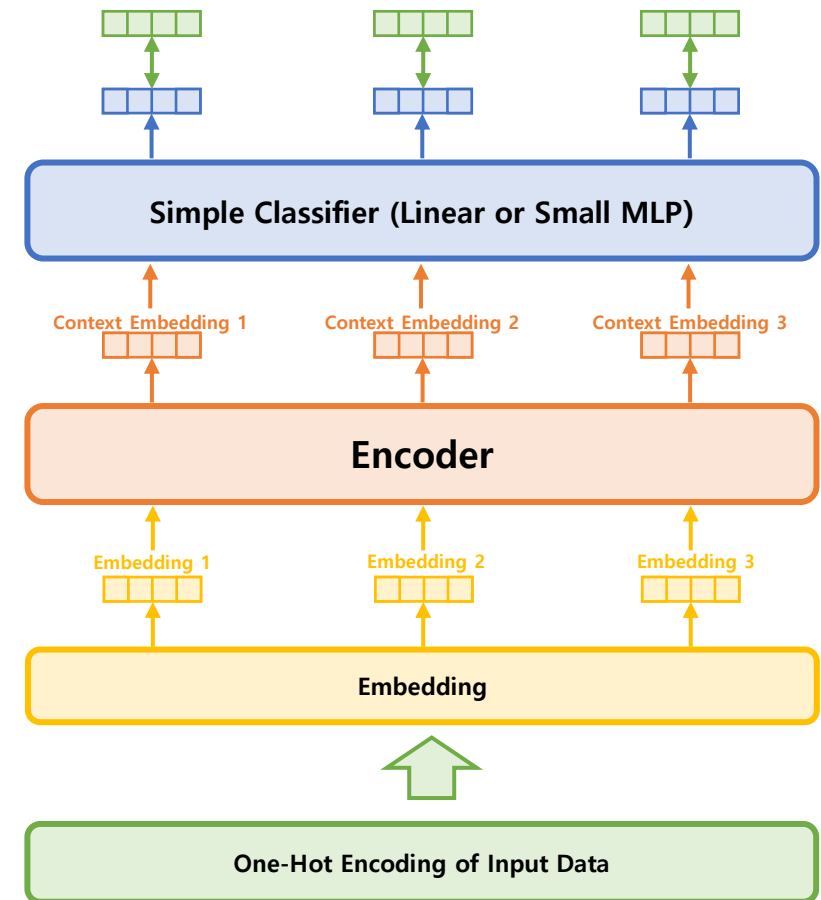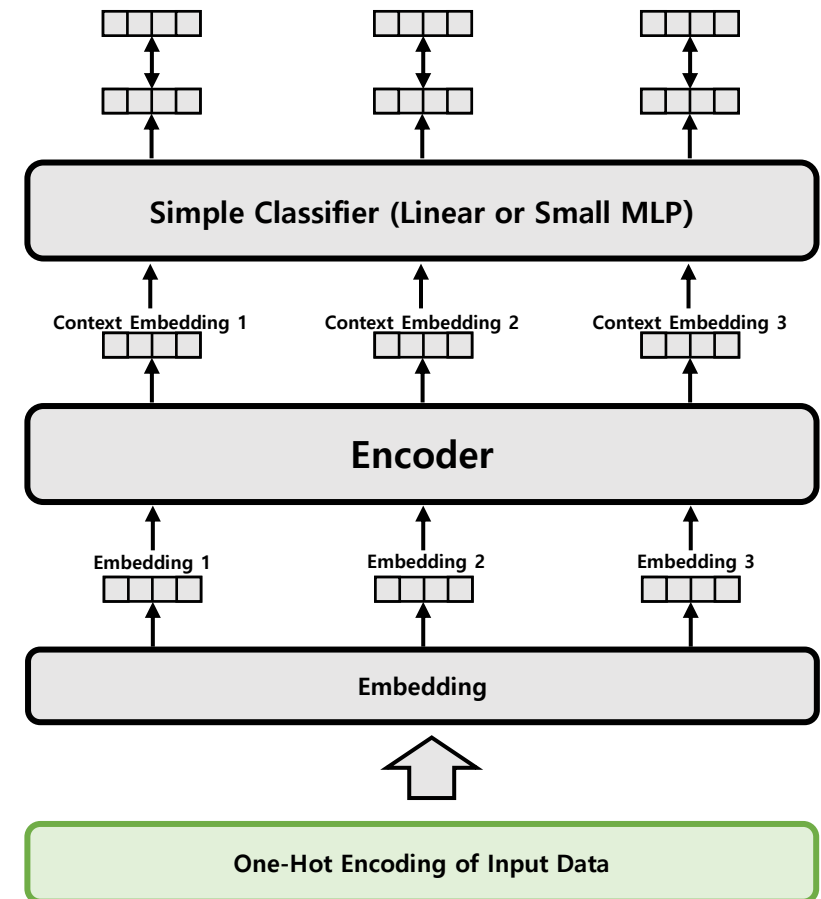
$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_{\theta}(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$



Context Embedding 1  Context Embedding 2  Context Embedding 3

Simple Classifier (Linear or Small MLP)

Encoder

Embedding 1  Embedding 2  Embedding 3

Embedding

One-Hot Encoding of Input Data

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$
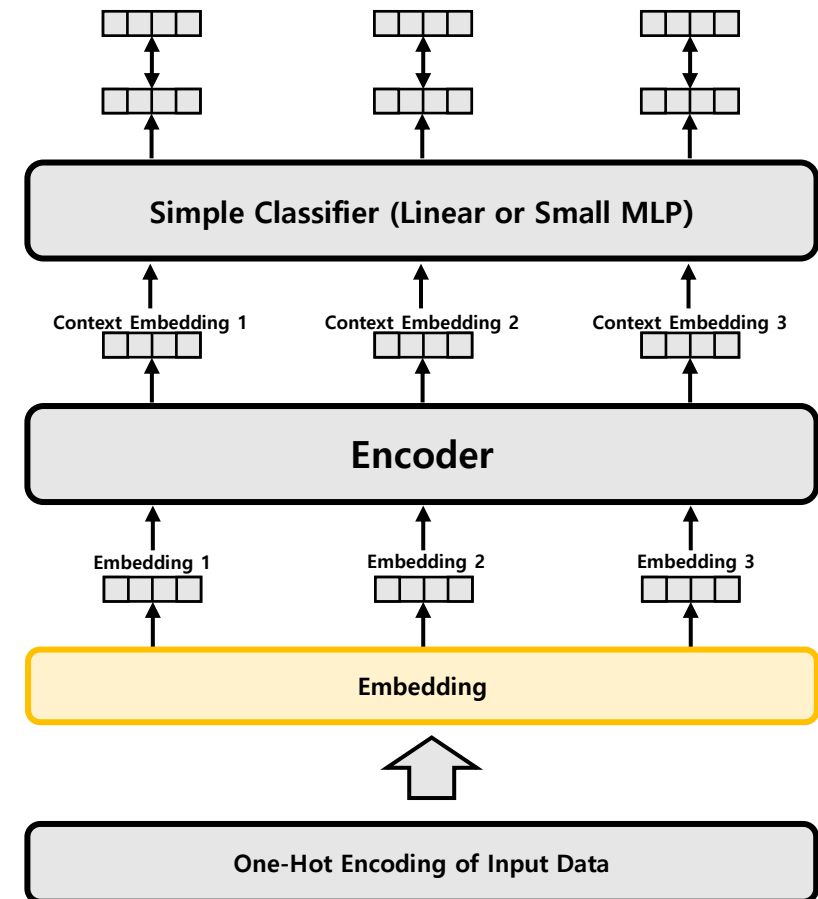
$Z: One - Hot\ Encoding$

$V: Embedding\ Matrix$

$X = VZ: Embedding$

$f_{\theta}(X): Language\ Model\ (Encoder)\ as\ Function$

$\theta: All\ Learnable\ Parameter\ in\ Language\ Model$

$y: Label$

$\delta: Perturbation$

**Simple Classifier (Linear or Small MLP)**

Context Embedding 1   Context Embedding 2   Context Embedding 3

**Encoder**

Embedding 1   Embedding 2   Embedding 3

**Embedding**

**One-Hot Encoding of Input Data**

# \<Projected Gradient Descent\>

$$\min_{\theta} \mathbb{E}_{(Z,\,y)\sim\mathcal{D}} \left[ \max_{\|\boldsymbol{\delta}\|\leq\boldsymbol{\varepsilon}} \boldsymbol{L(f_\theta(X+\delta),y)} \right]$$

$$\delta_{t+1} = \Pi_{\|\delta\|_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) \,/\, \big\|g(\delta_t)\big\|_F)$$

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \left||g(\delta_t)|\right|_F)$$

# <Projected Gradient Descent>

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\boldsymbol{Z}, \boldsymbol{y}) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} \boldsymbol{L}(\boldsymbol{f_{\theta}}(\boldsymbol{X} + \boldsymbol{\delta}), \boldsymbol{y}) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \left|\left|g(\delta_t)\right|\right|_F)$$

<Projected Gradient Descent>

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\boldsymbol{Z},\boldsymbol{y})\sim\mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} \boldsymbol{L}(\boldsymbol{f_\theta}(\boldsymbol{X}+\boldsymbol{\delta}),\boldsymbol{y}) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \left|\left|g(\delta_t)\right|\right|_F)$$

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z,\,y)\sim\mathcal{D}} \left[ \max_{\|\delta\| \le \varepsilon} L(f_\theta(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{\|\delta\|_F \le \varepsilon}(\delta_t + \alpha g(\delta_t) / \|g(\delta_t)\|_F)$$



Cat

Dog

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z,\,y)\sim\mathcal{D}} \left[ \max_{||\boldsymbol{\delta}|| \le \boldsymbol{\varepsilon}} L(f_{\theta}(X + \delta), y) \right]$$

$$\boldsymbol{\delta_{t+1}} = \boldsymbol{\Pi}_{||\boldsymbol{\delta}||_F \le \boldsymbol{\varepsilon}} (\boldsymbol{\delta_t} + \boldsymbol{\alpha g(\delta_t)} \,/\, ||\boldsymbol{g(\delta_t)}||_F)$$

# \<Projected Gradient Descent\>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$



Cat

Dog

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$



Cat

Dog

# <Projected Gradient Descent>

$$\min_\theta \mathbb{E}_{(Z,\,y)\sim\mathcal{D}} \left[ \max_{||\delta|| \le \varepsilon} L(f_\theta(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \le \varepsilon}(\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$
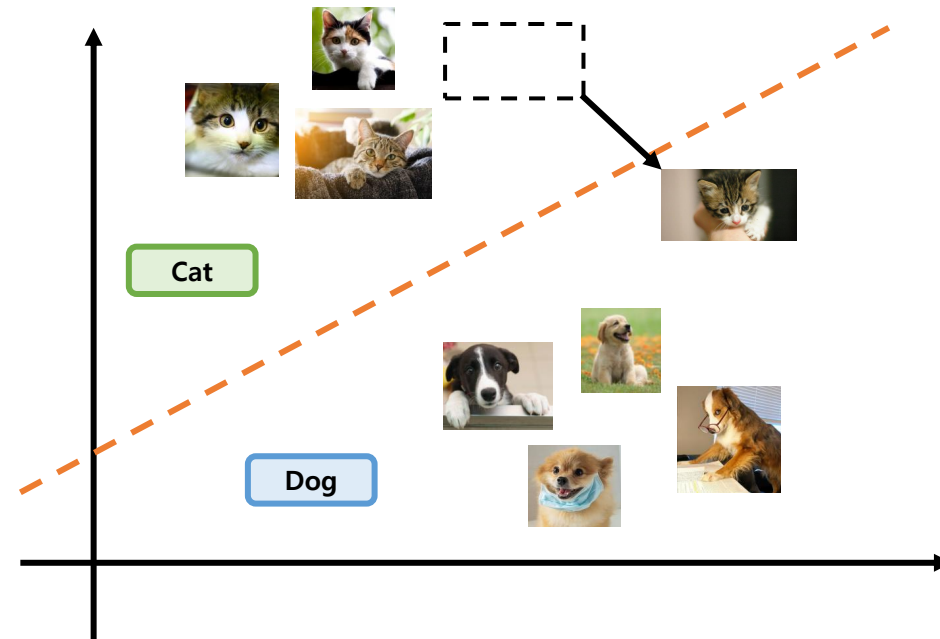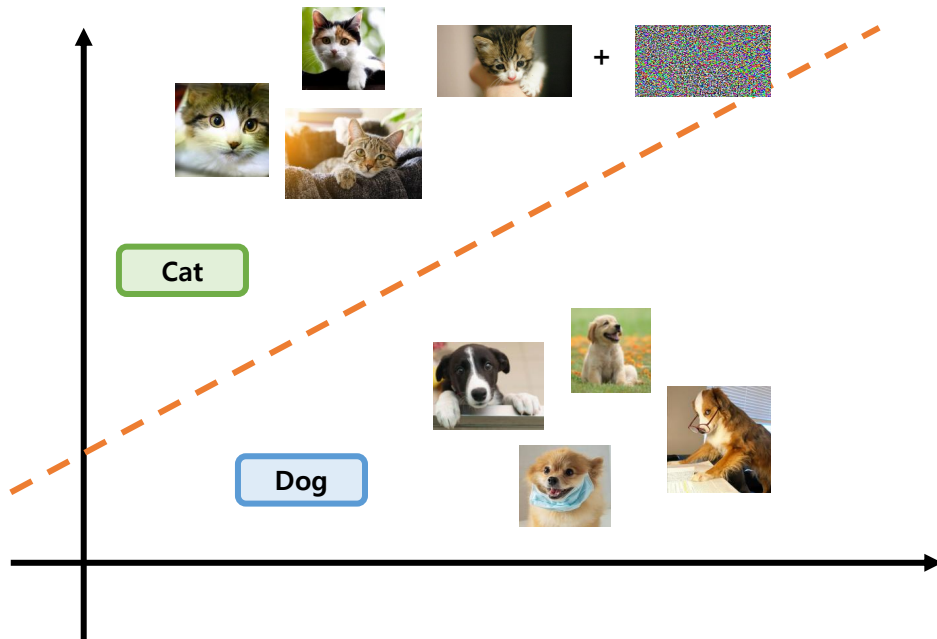
# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z,\, y)\sim \mathcal{D}} \left[ \max_{\|\delta\| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{\|\delta\|_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \|g(\delta_t)\|_F)$$

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{\|\delta\| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{\|\delta\|_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \|g(\delta_t)\|_F)$$

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z,\,y)\sim\mathcal{D}}\left[\max_{\|\delta\|\leq\varepsilon} L(f_{\theta}(X+\delta), y)\right]$$

$$\delta_{t+1} = \Pi_{\|\delta\|_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t)\,/\,\|g(\delta_t)\|_F)$$
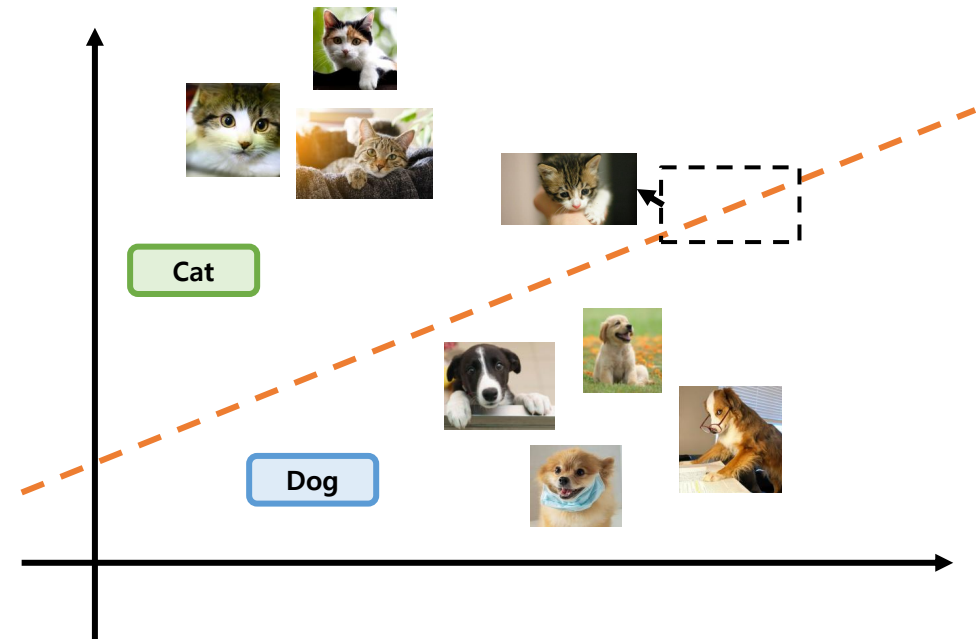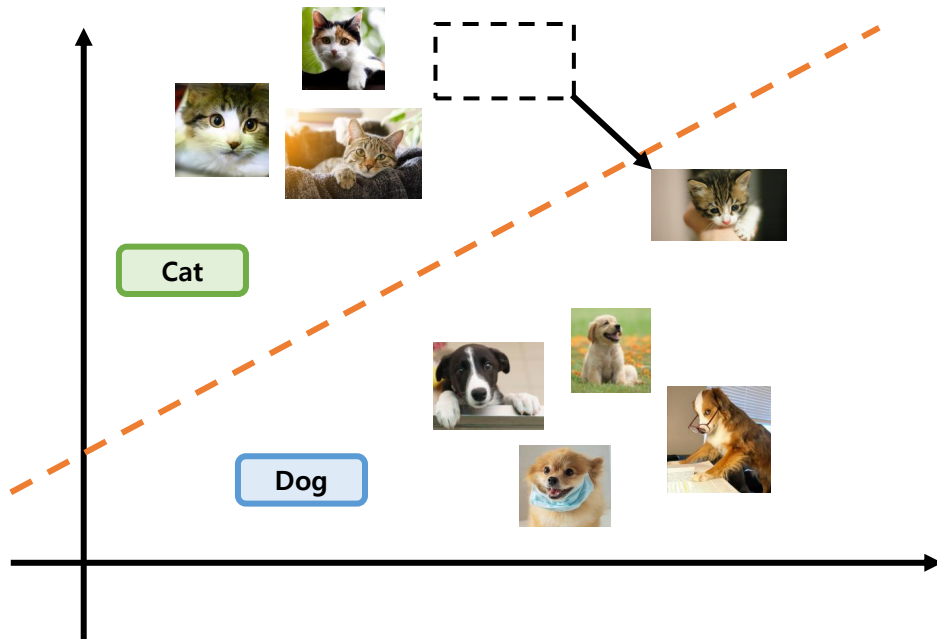
# \<Projected Gradient Descent\>

$$\min_{\theta} \mathbb{E}_{(Z,\,y)\sim\mathcal{D}} \left[ \max_{\|\delta\| \leq \varepsilon} L(f_\theta(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{\|\delta\|_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / \|g(\delta_t)\|_F)$$

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z,y)\sim\mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X+\delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$

# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{\|\delta\| \le \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{\|\delta\|_F \le \varepsilon}(\delta_t + \alpha g(\delta_t) / \|g(\delta_t)\|_F)$$
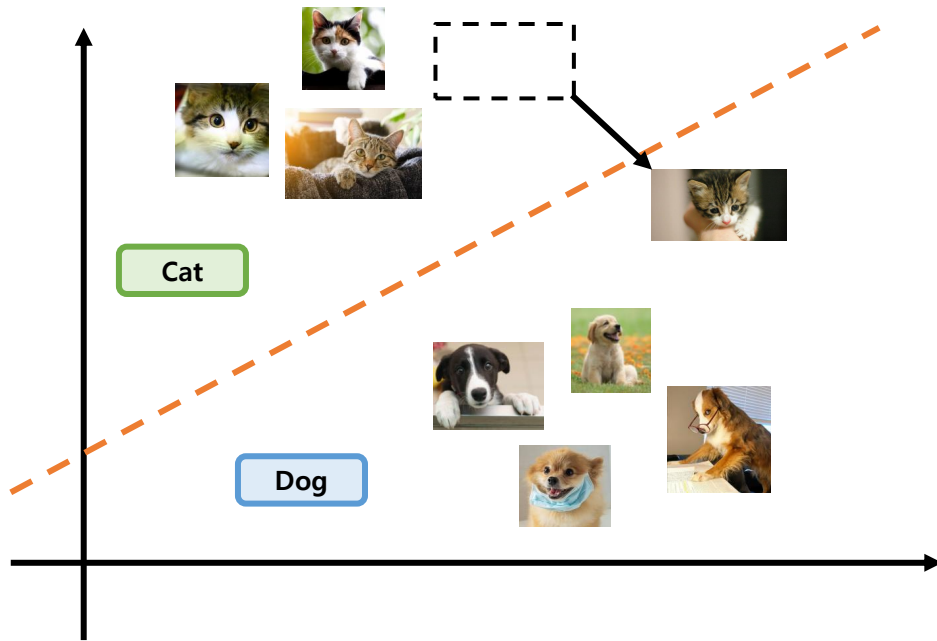
# <Projected Gradient Descent>

$$\min_{\theta} \mathbb{E}_{(Z, y) \sim \mathcal{D}} \left[ \max_{||\delta|| \leq \varepsilon} L(f_{\theta}(X + \delta), y) \right]$$

$$\delta_{t+1} = \Pi_{||\delta||_F \leq \varepsilon}(\delta_t + \alpha g(\delta_t) / ||g(\delta_t)||_F)$$
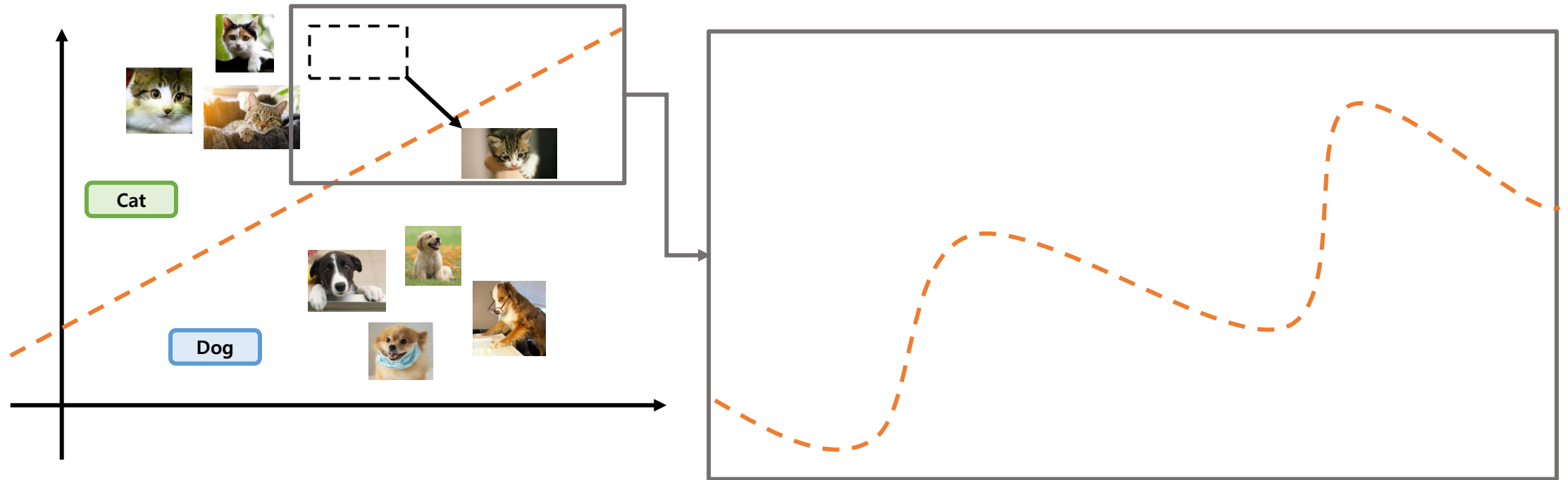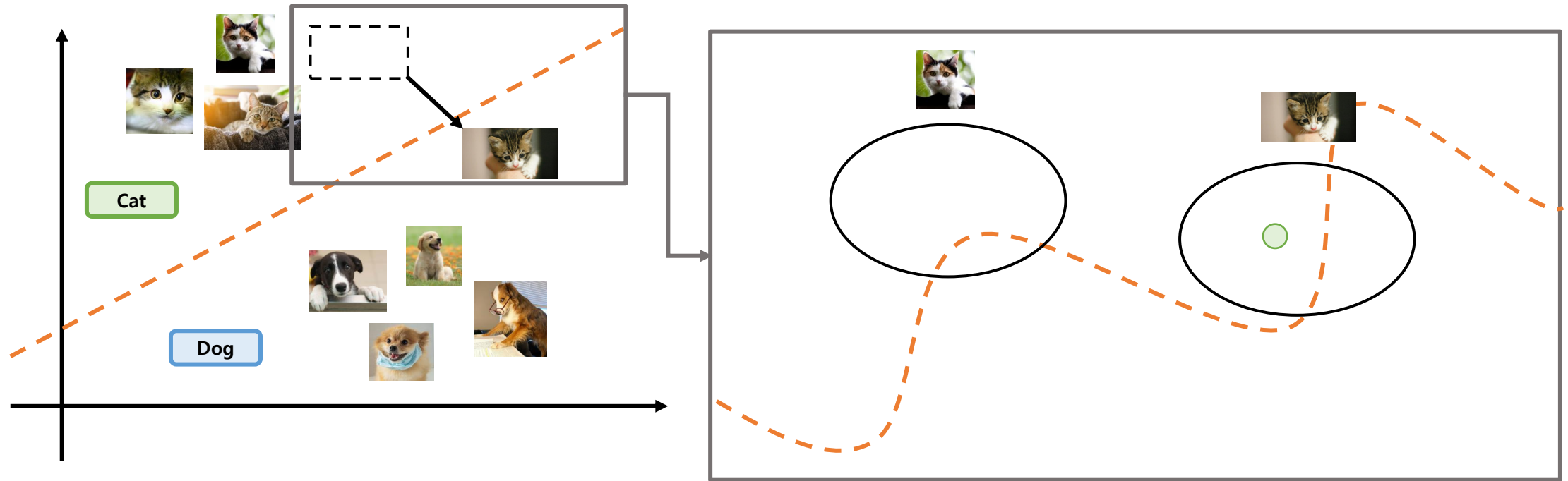
# &lt;FreeLB&gt;

| Algorithm | Perturbation | Model Parameter | Note |
|:---:|:---:|:---:|:---|
| **PGD** | Update K-Steps | Update 1-Step | • Update Perturbation K-Times and Update Model Parameter 1-Time |
| **FreeAT** | Update 1-Step | Update 1-Step | • Update Perturbation 1-Time and Update Model Parameter 1-Time |
| **YOPO** | Update K-Steps | Update 1-Step | • Accumulate Gradients for Model Parameters While Updating Perturbation K-Times and Update Model Parameters 1-Time Using Accumulated Gradients |

# \<FreeLB\>

---

**Algorithm 1** "Free" Large-Batch Adversarial Training (FreeLB-$K$)

---

**Require:** Training samples $X = \{(\boldsymbol{Z}, y)\}$, perturbation bound $\epsilon$, learning rate $\tau$, ascent steps $K$, ascent step size $\alpha$

1: Initialize $\boldsymbol{\theta}$
2: **for** epoch $= 1 \ldots N_{ep}$ **do**
3:      **for** minibatch $B \subset X$ **do**
4:         $\boldsymbol{\delta}_0 \leftarrow \frac{1}{\sqrt{N_\delta}} U(-\epsilon, \epsilon)$
5:         $\boldsymbol{g}_0 \leftarrow 0$
6:         **for** $t = 1 \ldots K$ **do**
7:             Accumulate gradient of parameters $\theta$
8:             $\boldsymbol{g}_t \leftarrow \boldsymbol{g}_{t-1} + \frac{1}{K} \mathbb{E}_{(\boldsymbol{Z}, y) \in B} [\nabla_{\boldsymbol{\theta}} L(f_{\boldsymbol{\theta}}(\boldsymbol{X} + \boldsymbol{\delta}_{t-1}), y)]$
9:             Update the perturbation $\delta$ via gradient ascend
10:            $\boldsymbol{g}_{adv} \leftarrow \nabla_{\boldsymbol{\delta}} L(f_{\boldsymbol{\theta}}(\boldsymbol{X} + \boldsymbol{\delta}_{t-1}), y)$
11:            $\boldsymbol{\delta}_t \leftarrow \Pi_{\|\boldsymbol{\delta}\|_F \leq \epsilon}(\boldsymbol{\delta}_{t-1} + \alpha \cdot \boldsymbol{g}_{adv} / \|\boldsymbol{g}_{adv}\|_F)$
12:         **end for**
13:         $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \tau \boldsymbol{g}_K$
14:      **end for**
15: **end for**

---

# \<FreeLB\>

---

**Algorithm 1** "Free" Large-Batch Adversarial Training (FreeLB-$K$)

---

**Require:** Training samples $X = \{(\boldsymbol{Z}, y)\}$, perturbation bound $\epsilon$, learning rate $\tau$, ascent steps $K$, ascent step size $\alpha$

1: Initialize $\boldsymbol{\theta}$
2: **for** epoch $= 1 \ldots N_{ep}$ **do**
3:     **for** minibatch $B \subset X$ **do**
4:         $\boldsymbol{\delta}_0 \leftarrow \frac{1}{\sqrt{N_\delta}} U(-\epsilon, \epsilon)$
5:         $\boldsymbol{g}_0 \leftarrow 0$
6:         **for** $t = 1 \ldots K$ **do**
7:             Accumulate gradient of parameters $\theta$
8:             $\boldsymbol{g}_t \leftarrow \boldsymbol{g}_{t-1} + \frac{1}{K} \mathbb{E}_{(\boldsymbol{Z}, y) \in B} [\nabla_{\boldsymbol{\theta}} L(f_{\boldsymbol{\theta}}(\boldsymbol{X} + \boldsymbol{\delta}_{t-1}), y)]$
9:             Update the perturbation $\delta$ via gradient ascend
10:            $\boldsymbol{g}_{adv} \leftarrow \nabla_{\boldsymbol{\delta}} L(f_{\boldsymbol{\theta}}(\boldsymbol{X} + \boldsymbol{\delta}_{t-1}), y)$
11:            $\boldsymbol{\delta}_t \leftarrow \Pi_{\|\boldsymbol{\delta}\|_F \leq \epsilon} (\boldsymbol{\delta}_{t-1} + \alpha \cdot \boldsymbol{g}_{adv} / \|\boldsymbol{g}_{adv}\|_F)$
12:         **end for**
13:         $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \tau \boldsymbol{g}_K$
14:     **end for**
15: **end for**

---

# Experiments

- GLUE Benchmark

- Comparing the Robustness

# \<GLUE Benchmark\>

| Method | MNLI | QNLI | QQP | RTE | SST-2 | MRPC | CoLA | STS-B |
|---|---|---|---|---|---|---|---|---|
| | (Acc) | (Acc) | (Acc) | (Acc) | (Acc) | (Acc) | (Mcc) | (Pearson) |
| Reported | 90.20 | 94.70 | 92.20 | 86.60 | 96.40 | 90.90 | 68.00 | 92.40 |
| ReImp | - | - | - | 85.61 (1.7) | 96.56 (.3) | 90.69 (.5) | 67.57 (1.3) | 92.20 (.2) |
| PGD | 90.53 (.2) | 94.87 (.2) | 92.49 (.07) | 87.41 (.9) | 96.44 (.1) | 90.93 (.2) | 69.67 (1.2) | 92.43 (7.) |
| FreeAT | 90.02 (.2) | 94.66 (.2) | 92.48 (.08) | 86.69 (15.) | 96.10 (.2) | 90.69 (.4) | 68.80 (1.3) | 92.40 (.3) |
| FreeLB | **90.61** (.1) | **94.98** (.2) | **92.60** (03) | **88.13** (1.2) | **96.79** (.2) | **91.41** (.7) | **71.12** (.9) | **92.67** (.08) |

**\<Results (Median and Variance) on the dev sets of GLUE based on the RoBERTa-Large Model\>**

**Adversarial Training for NLU**

- GLUE Benchmark

# \<GLUE Benchmark\>

| Method | MNLI (Acc) | QNLI (Acc) | QQP (Acc) | RTE (Acc) | SST-2 (Acc) | MRPC (Acc) | CoLA (Mcc) | STS-B (Pearson) |
|---|---|---|---|---|---|---|---|---|
| Reported | 90.20 | 94.70 | 92.20 | 86.60 | 96.40 | 90.90 | 68.00 | 92.40 |
| ReImp | - | - | - | 85.61 (1.7) | 96.56 (.3) | 90.69 (.5) | 67.57 (1.3) | 92.20 (.2) |
| PGD | 90.53 (.2) | 94.87 (.2) | 92.49 (.07) | 87.41 (.9) | 96.44 (.1) | 90.93 (.2) | 69.67 (1.2) | 92.43 (7.) |
| FreeAT | 90.02 (.2) | 94.66 (.2) | 92.48 (.08) | 86.69 (15.) | 96.10 (.2) | 90.69 (.4) | 68.80 (1.3) | 92.40 (.3) |
| FreeLB | **90.61** (.1) | **94.98** (.2) | **92.60** (03) | **88.13** (1.2) | **96.79** (.2) | **91.41** (.7) | **71.12** (.9) | **92.67** (.08) |

**Results (Median and Variance) on the dev sets of GLUE based on the RoBERTa-Large Model**

# <GLUE Benchmark>

| Method | MNLI | QNLI | QQP | RTE | SST-2 | MRPC | CoLA | STS-B |
|---|---|---|---|---|---|---|---|---|
| | (Acc) | (Acc) | (Acc) | (Acc) | (Acc) | (Acc) | (Mcc) | (Pearson) |
| Reported | 90.20 | 94.70 | 92.20 | 86.60 | 96.40 | 90.90 | 68.00 | 92.40 |
| ReImp | - | - | - | 85.61 (1.7) | 96.56 (.3) | 90.69 (.5) | 67.57 (1.3) | 92.20 (.2) |
| PGD | 90.53 (.2) | 94.87 (.2) | 92.49 (.07) | 87.41 (.9) | 96.44 (.1) | 90.93 (.2) | 69.67 (1.2) | 92.43 (7.) |
| FreeAT | 90.02 (.2) | 94.66 (.2) | 92.48 (.08) | 86.69 (15.) | 96.10 (.2) | 90.69 (.4) | 68.80 (1.3) | 92.40 (.3) |
| **FreeLB** | **90.61** (.1) | **94.98** (.2) | **92.60** (03) | **88.13** (1.2) | **96.79** (.2) | **91.41** (.7) | **71.12** (.9) | **92.67** (.08) |

**Results (Median and Variance) on the dev sets of GLUE based on the RoBERTa-Large Model**

# Adversarial Training for NLU
**- GLUE Benchmark**

## \<GLUE Benchmark\>

| Model | Score | CoLA 8.5k | SST-2 67k | MRPC 3.7k | STS-B 7k | QQP 364k | MNLI-m/mm 393k | QNLI 108k | RTE 2.5k | WNLI 634 | AX |
|---|---|---|---|---|---|---|---|---|---|---|---|
| BERT-base | 78.3 | 52.1 | 93.5 | 88.9/88.4 | 87.1/85.8 | 71.2/89.2 | 94.6/83.4 | 90.5 | 66.4 | 65.1 | 64.2 |
| FreeLB-BERT | 79.4 | 54.5 | 93.6 | 88.1/83.5 | 87.7/86.7 | 72.7/89.6 | 85.7/84.6 | 91.8 | 70.1 | 65.1 | 36.9 |
| MT-DNN | 87.6 | **68.4** | 96.5 | 92.7/90.3 | 91.1/90.7 | 73.7/89.9 | 87.9/87.4 | 96.0 | 86.3 | 89.0 | 42.8 |
| XLNet-Large | 88.4 | 67.8 | **96.8** | 93.0/90.7 | 91.6/91.1 | 74.2/90.3 | 90.2/89.8 | 98.6 | 86.3 | **90.4** | 47.5 |
| RoBERTa | 88.5 | 67.8 | 96.7 | 92.3/89.8 | 92.2/91.9 | 74.3/90.2 | 90.8/90.2 | **98.9** | 88.2 | 89.0 | 47.7 |
| FreeLB-RoB | **88.8** | 68.0 | **96.8** | **93.1/90.8** | **92.4/92.2** | **74.8/90.3** | **91.1/90.7** | 98.8 | **88.7** | 89.0 | **50.1** |
| Human | 87.1 | 66.4 | 97.8 | 86.3/80.8 | 92.7/92.6 | 59.5/80.4 | 92.0/92.8 | 91.2 | 93.6 | 95.9 | - |

**Results on GLUE from the Evaluation Server, as of Sep 25, 2019**

# Adversarial Training for NLU
## - GLUE Benchmark

## <GLUE Benchmark>

| Model | Score | CoLA 8.5k | SST-2 67k | MRPC 3.7k | STS-B 7k | QQP 364k | MNLI-m/mm 393k | QNLI 108k | RTE 2.5k | WNLI 634 | AX |
|---|---|---|---|---|---|---|---|---|---|---|---|
| BERT-base | 78.3 | 52.1 | 93.5 | 88.9/88.4 | 87.1/85.8 | 71.2/89.2 | 94.6/83.4 | 90.5 | 66.4 | 65.1 | 64.2 |
| FreeLB-BERT | 79.4 | 54.5 | 93.6 | 88.1/83.5 | 87.7/86.7 | 72.7/89.6 | 85.7/84.6 | 91.8 | 70.1 | 65.1 | 36.9 |
| MT-DNN | 87.6 | **68.4** | 96.5 | 92.7/90.3 | 91.1/90.7 | 73.7/89.9 | 87.9/87.4 | 96.0 | 86.3 | 89.0 | 42.8 |
| XLNet-Large | 88.4 | 67.8 | **96.8** | 93.0/90.7 | 91.6/91.1 | 74.2/90.3 | 90.2/89.8 | 98.6 | 86.3 | **90.4** | 47.5 |
| RoBERTa | 88.5 | 67.8 | 96.7 | 92.3/89.8 | 92.2/91.9 | 74.3/90.2 | 90.8/90.2 | **98.9** | 88.2 | 89.0 | 47.7 |
| FreeLB-RoB | **88.8** | 68.0 | **96.8** | **93.1/90.8** | **92.4/92.2** | **74.8/90.3** | **91.1/90.7** | 98.8 | **88.7** | 89.0 | **50.1** |
| Human | 87.1 | 66.4 | 97.8 | 86.3/80.8 | 92.7/92.6 | 59.5/80.4 | 92.0/92.8 | 91.2 | 93.6 | 95.9 | - |

**Results on GLUE from the Evaluation Server, as of Sep 25, 2019**

# Adversarial Training for NLU
- GLUE Benchmark

## \<GLUE Benchmark\>

| Model | Score | CoLA 8.5k | SST-2 67k | MRPC 3.7k | STS-B 7k | QQP 364k | MNLI-m/mm 393k | QNLI 108k | RTE 2.5k | WNLI 634 | AX |
|-------|-------|-----------|-----------|-----------|----------|----------|----------------|-----------|----------|----------|------|
| **BERT-base** | 78.3 | 52.1 | 93.5 | 88.9/88.4 | 87.1/85.8 | 71.2/89.2 | 94.6/83.4 | 90.5 | 66.4 | 65.1 | 64.2 |
| **FreeLB-BERT** | 79.4 | 54.5 | 93.6 | 88.1/83.5 | 87.7/86.7 | 72.7/89.6 | 85.7/84.6 | 91.8 | 70.1 | 65.1 | 36.9 |
| **MT-DNN** | 87.6 | **68.4** | 96.5 | 92.7/90.3 | 91.1/90.7 | 73.7/89.9 | 87.9/87.4 | 96.0 | 86.3 | 89.0 | 42.8 |
| **XLNet-Large** | 88.4 | 67.8 | **96.8** | 93.0/90.7 | 91.6/91.1 | 74.2/90.3 | 90.2/89.8 | 98.6 | 86.3 | **90.4** | 47.5 |
| **RoBERTa** | 88.5 | 67.8 | 96.7 | 92.3/89.8 | 92.2/91.9 | 74.3/90.2 | 90.8/90.2 | **98.9** | 88.2 | 89.0 | 47.7 |
| **FreeLB-RoB** | **88.8** | 68.0 | **96.8** | **93.1/90.8** | **92.4/92.2** | **74.8/90.3** | **91.1/90.7** | 98.8 | **88.7** | 89.0 | **50.1** |
| **Human** | 87.1 | 66.4 | 97.8 | 86.3/80.8 | 92.7/92.6 | 59.5/80.4 | 92.0/92.8 | 91.2 | 93.6 | 95.9 | - |

**Results on GLUE from the Evaluation Server, as of Sep 25, 2019**

# Adversarial Training for NLU
- GLUE Benchmark

## <GLUE Benchmark>

| Model | Score | CoLA 8.5k | SST-2 67k | MRPC 3.7k | STS-B 7k | QQP 364k | MNLI-m/mm 393k | QNLI 108k | RTE 2.5k | WNLI 634 | AX |
|---|---|---|---|---|---|---|---|---|---|---|---|
| BERT-base | 78.3 | 52.1 | 93.5 | 88.9/88.4 | 87.1/85.8 | 71.2/89.2 | 94.6/83.4 | 90.5 | 66.4 | 65.1 | 64.2 |
| FreeLB-BERT | 79.4 | 54.5 | 93.6 | 88.1/83.5 | 87.7/86.7 | 72.7/89.6 | 85.7/84.6 | 91.8 | 70.1 | 65.1 | 36.9 |
| MT-DNN | 87.6 | **68.4** | 96.5 | 92.7/90.3 | 91.1/90.7 | 73.7/89.9 | 87.9/87.4 | 96.0 | 86.3 | 89.0 | 42.8 |
| XLNet-Large | 88.4 | 67.8 | **96.8** | 93.0/90.7 | 91.6/91.1 | 74.2/90.3 | 90.2/89.8 | 98.6 | 86.3 | **90.4** | 47.5 |
| RoBERTa | 88.5 | 67.8 | 96.7 | 92.3/89.8 | 92.2/91.9 | 74.3/90.2 | 90.8/90.2 | **98.9** | 88.2 | 89.0 | 47.7 |
| FreeLB-RoB | **88.8** | 68.0 | **96.8** | **93.1/90.8** | **92.4/92.2** | **74.8/90.3** | **91.1/90.7** | 98.8 | **88.7** | 89.0 | **50.1** |
| Human | 87.1 | 66.4 | 97.8 | 86.3/80.8 | 92.7/92.6 | 59.5/80.4 | 92.0/92.8 | 91.2 | 93.6 | 95.9 | - |

## Results on GLUE from the Evaluation Server, as of Sep 25, 2019

# <Comparing the Robustness>

| Methods | RTE | | | CoLA | | | MRPC | | |
|---|---|---|---|---|---|---|---|---|---|
| | **M-Inc** $(10^{-4})$ | **M-Inc (R)** $(10^{-4})$ | **N-Loss** $(10^{-4})$ | **M-Inc** $(10^{-4})$ | **M-Inc (R)** $(10^{-4})$ | **N-Loss** $(10^{-4})$ | **M-Inc** $(10^{-3})$ | **M-Inc (R)** $(10^{-3})$ | **N-Loss** $(10^{-3})$ |
| **Vanilla** | 5.1 | 5.3 | 4.5 | 6.1 | 5.7 | 5.2 | 10.2 | 10.2 | 1.9 |
| **PGD** | 4.7 | 4.9 | 6.2 | 128.2 | 130.1 | 436.1 | 5.7 | 5.7 | 5.4 |
| **FreeLB** | 3.0 | 2.6 | 4.1 | 1.4 | 1.3 | 7.2 | 3.6 | 3.6 | 2.7 |

**Median of the Maximum Increase in Loss in the Vicinity of the Dev Set Samples for RoBERTa-Large Model Finetuned with Different Methods**

$$\Delta L_{max}(X, \epsilon) = \max_{||\delta|| \leq \epsilon} L(f_\theta(X + \delta), y) - L(f_\theta(X), y)$$

| | |
|---|---|
| M-Inc | FreeLB |
| M-Inc (R) | PGD |
| N-Loss | Clean Sample |

# <Comparing the Robustness>

| Methods | RTE | | | CoLA | | | MRPC | | |
|---|---|---|---|---|---|---|---|---|---|
| | M-Inc $(10^{-4})$ | M-Inc (R) $(10^{-4})$ | N-Loss $(10^{-4})$ | M-Inc $(10^{-4})$ | M-Inc (R) $(10^{-4})$ | N-Loss $(10^{-4})$ | M-Inc $(10^{-3})$ | M-Inc (R) $(10^{-3})$ | N-Loss $(10^{-3})$ |
| Vanilla | 5.1 | 5.3 | 4.5 | 6.1 | 5.7 | 5.2 | 10.2 | 10.2 | 1.9 |
| PGD | 4.7 | 4.9 | 6.2 | 128.2 | 130.1 | 436.1 | 5.7 | 5.7 | 5.4 |
| FreeLB | 3.0 | 2.6 | 4.1 | 1.4 | 1.3 | 7.2 | 3.6 | 3.6 | 2.7 |

**Median of the Maximum Increase in Loss in the Vicinity of the Dev Set Samples for RoBERTa-Large Model Finetuned with Different Methods**

$$\Delta L_{max}(X, \epsilon) = \max_{||\delta|| \leq \epsilon} L(f_\theta(X + \delta), y) - L(f_\theta(X), y)$$

| | |
|---|---|
| M-Inc | FreeLB |
| M-Inc (R) | PGD |
| N-Loss | Clean Sample |

# Conclusion

# <Conclusion>

- Proposed a novel adversarial training algorithm, FreeLB, that promotes higher invariance in the embedding space

- Applied FreeLB to Transformer-based models for natural language understanding and achieved new state-of-the-art on GLUE benchmark

- FreeLB resulted in both higher robustness in the embedding space than natural training and better generalization ability

# Any Questions?

# Thank You