

The 38<sup>th</sup> ACM/SIGAPP Symposium On Applied Computing (Tallinn Estonia, March 27 – March 31, 2023)

---

# Security Verification Software Platform of Data-efficient Image Transformer Based on Fast Gradient Sign Method

In-pyo Hong, Gyu-ho Choi, Pan-Koo Kim and Chang Choi\*

Presenter : In-pyo Hong (Dept. of Computer Engineering, Gachon University, Republic of Korea)

# • Introduction

---

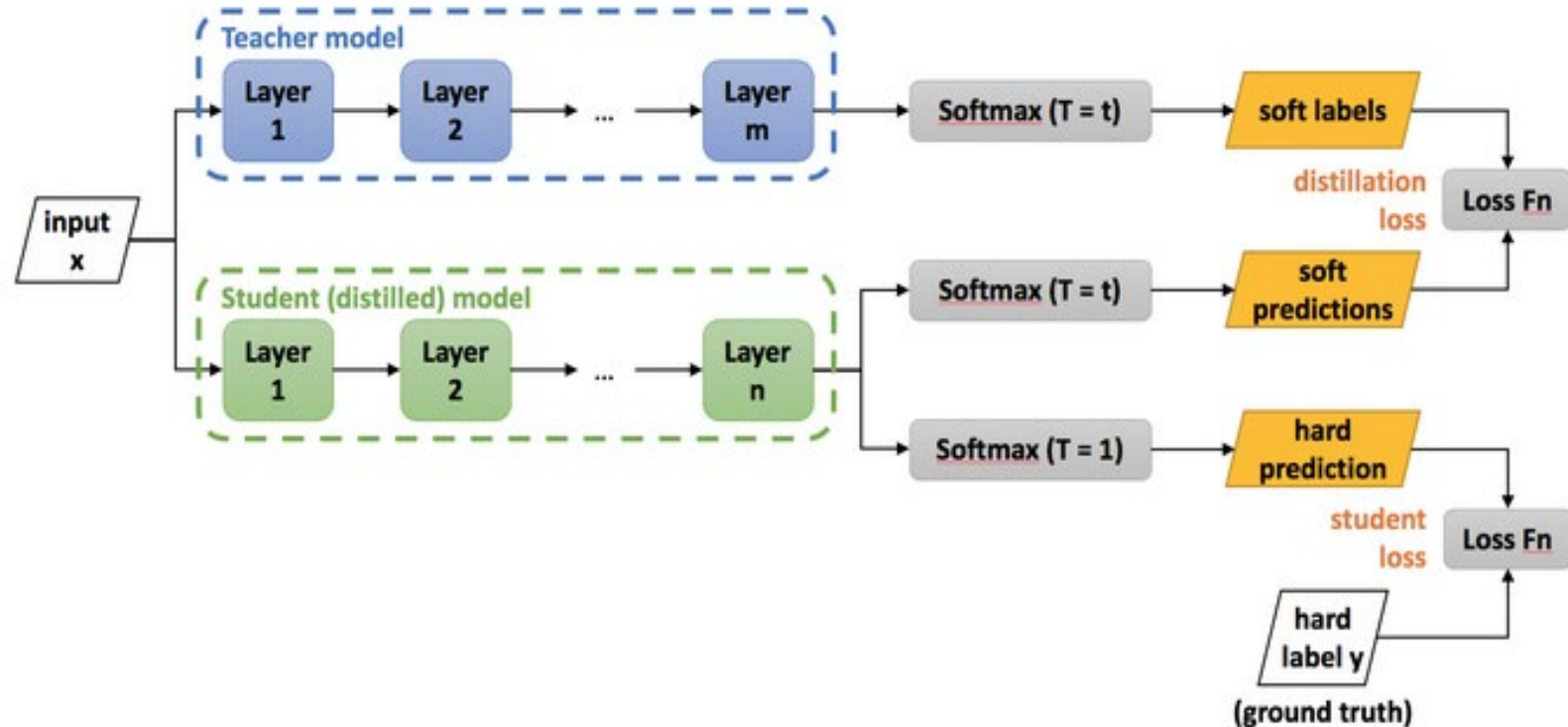
- **Knowledge distillation** in AI has been actively conducted.
- Representative Image transformer model of **knowledge distillation**  
→ **DeiT (Data-efficient Image Transformer)**
- DeiT has not been verified as safe against the **AI security threat**  
(especially **adversarial attack**)



Research Purpose

- *DeiT's adversarial attack vulnerability analysis*

# • Background #knowledge Distillation

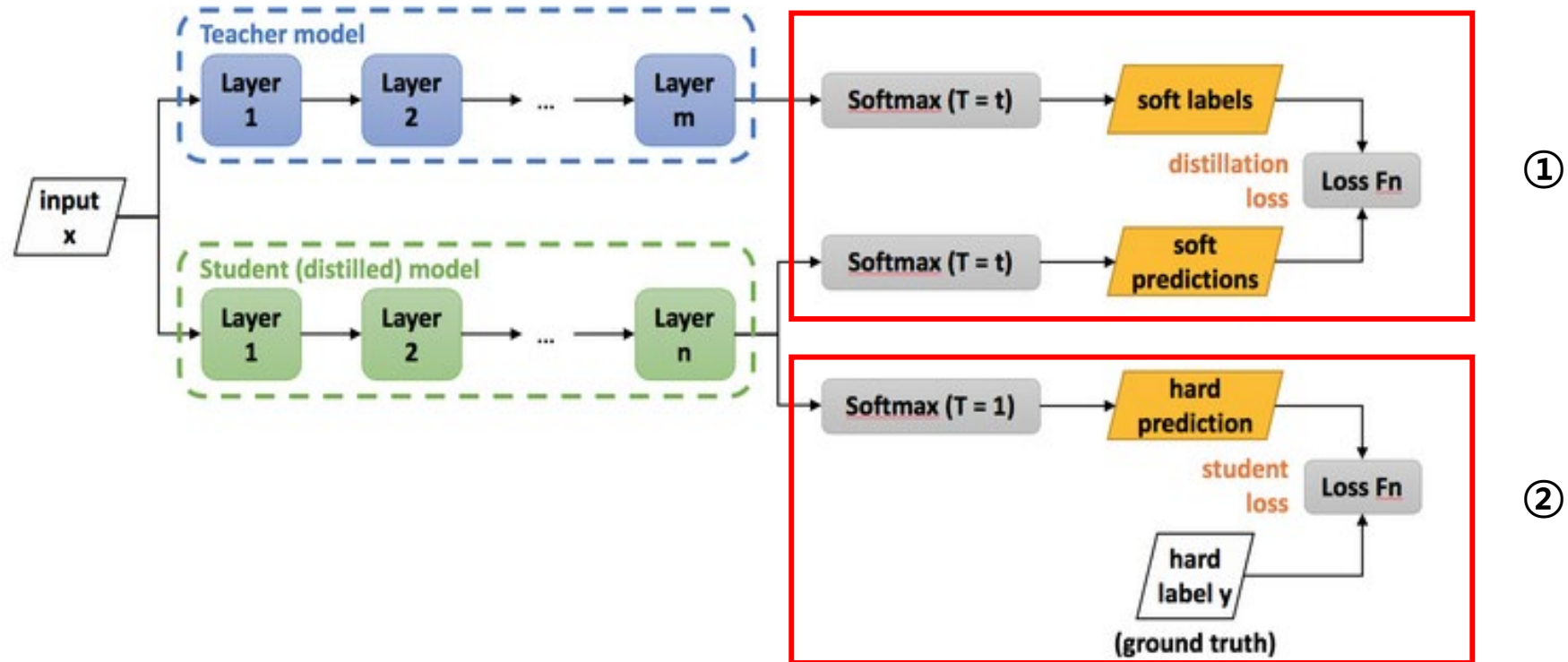


<Fig1> Knowledge Distillation Structure

Concept of knowledge distillation : Using **2 training Models**

Purpose of knowledge distillation : Maximizing the performance of the **Student Model**

# • Background #knowledge Distillation



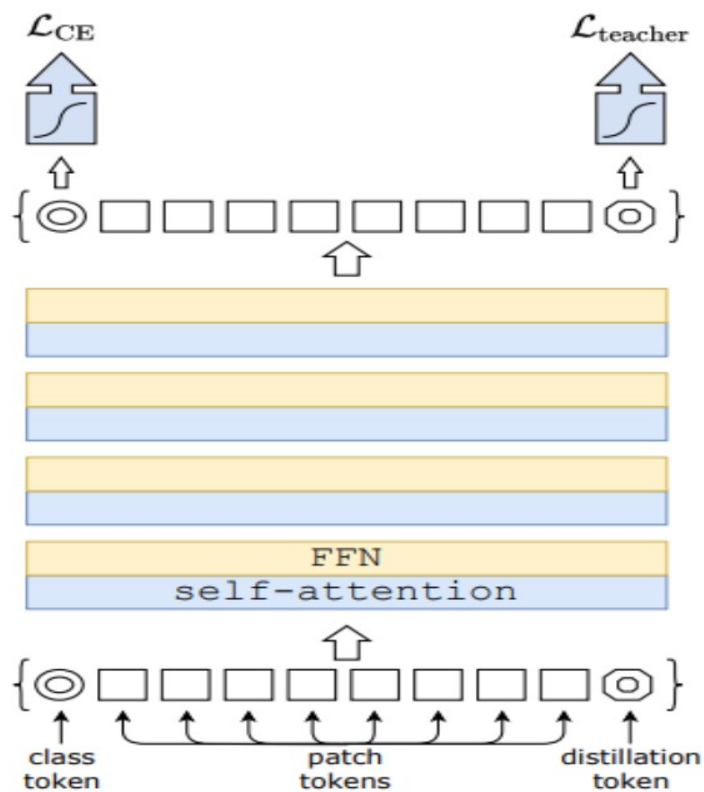
<Fig1> Knowledge Distillation Structure

Learning methods of the **Student Model** : **Twice Training** (① ②)

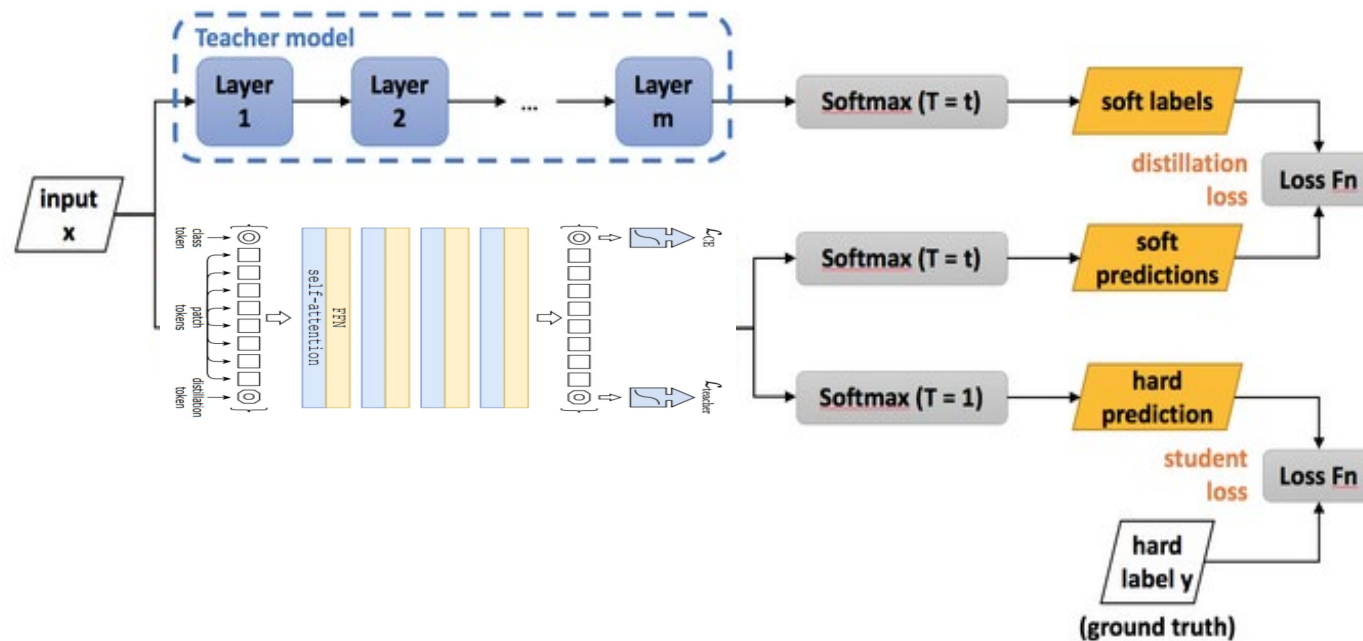
① : Training using the **soft labeled output of the teacher model** → Learning **the weights of the teacher model**

② : Basic supervised learning with ground truth

- Background #DeiT

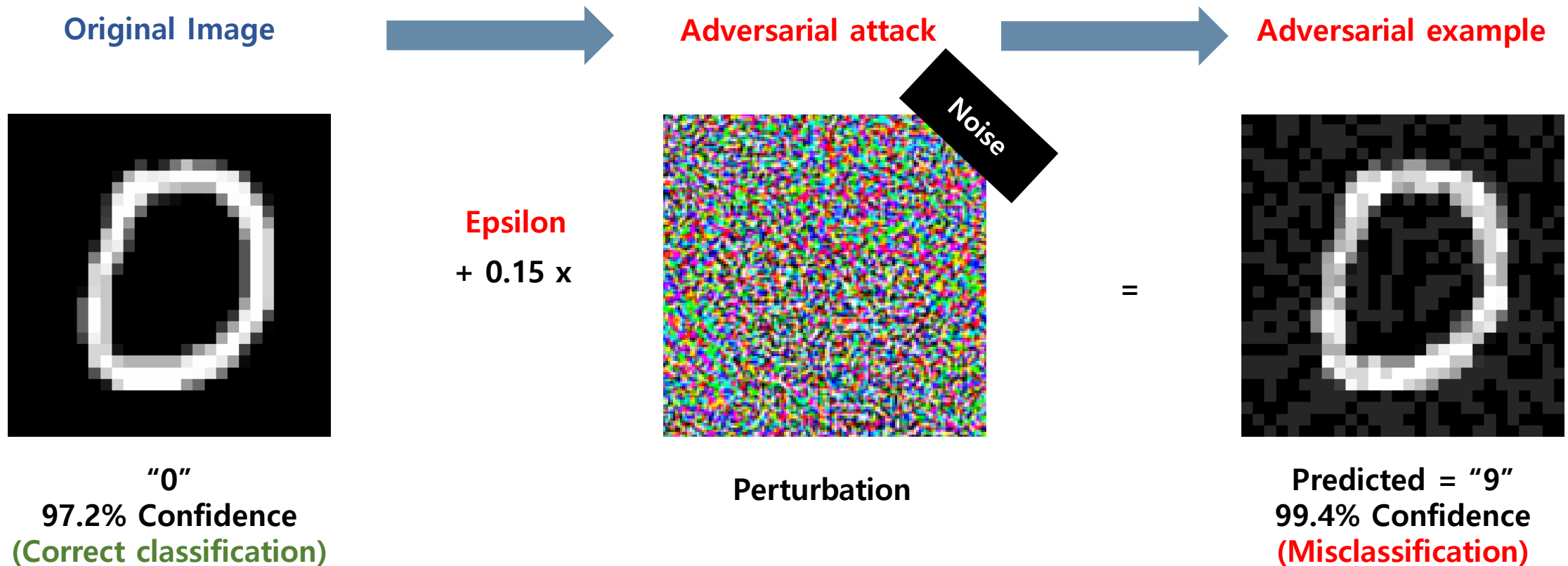


<Fig2> DeiT Model Structure



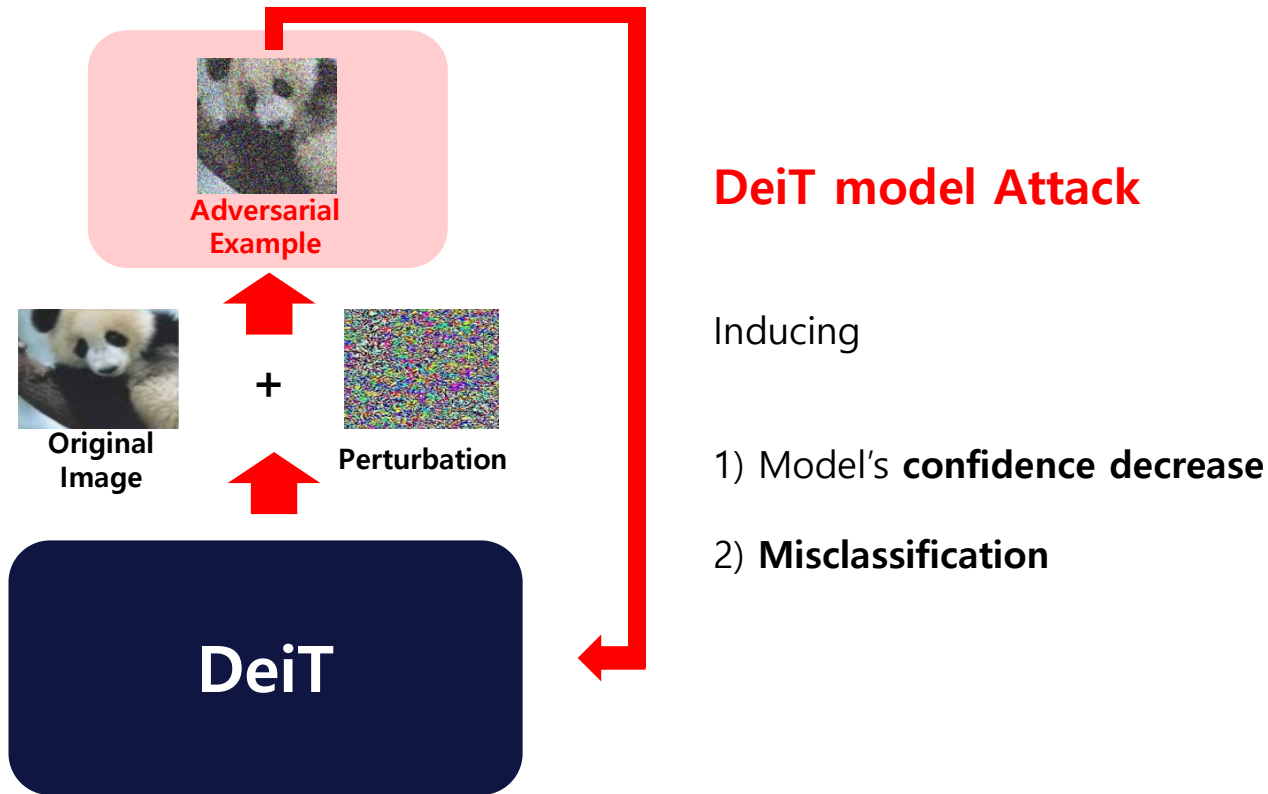
<Fig3> DeiT applied knowledge distillation structure

- Background #Adversarial Attack



<Fig4> Adversarial Attack Example

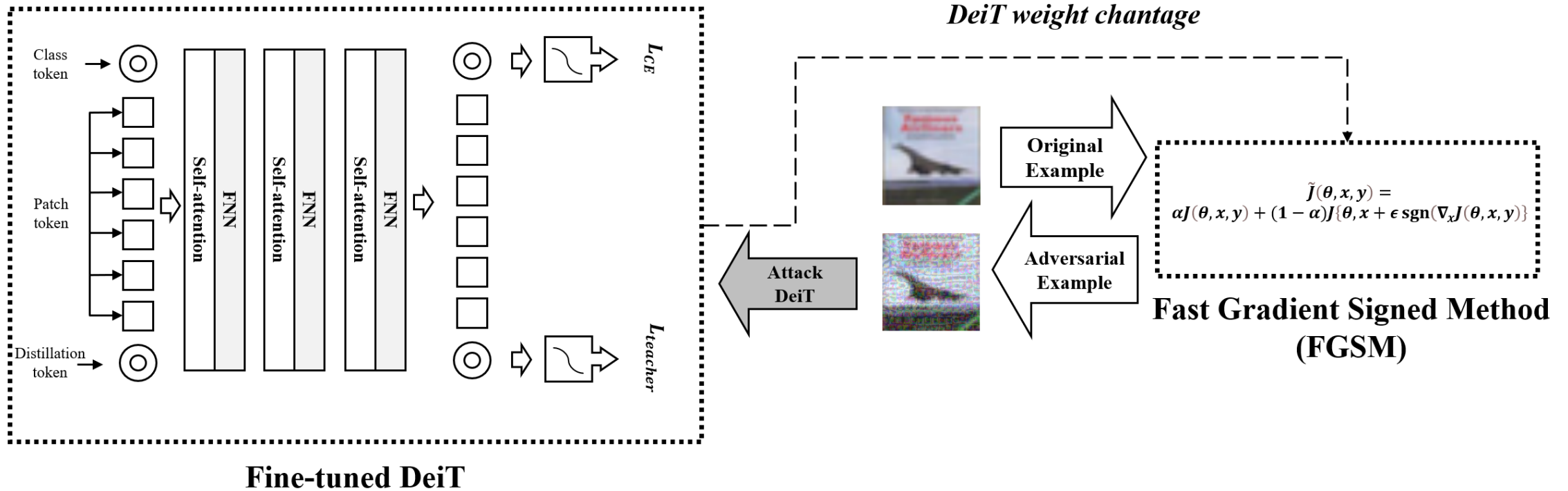
# • Background #Adversarial Attack



<Fig5> Adversarial Attack Process

- The **most dangerous attack** in the field of **AI security**
- **Typical adversarial attack technique**
  - 1) **FGSM**
  - 2) **PGD**
  - 3) **C&W Attack**

# • Method



<Fig6> Experiment Overview

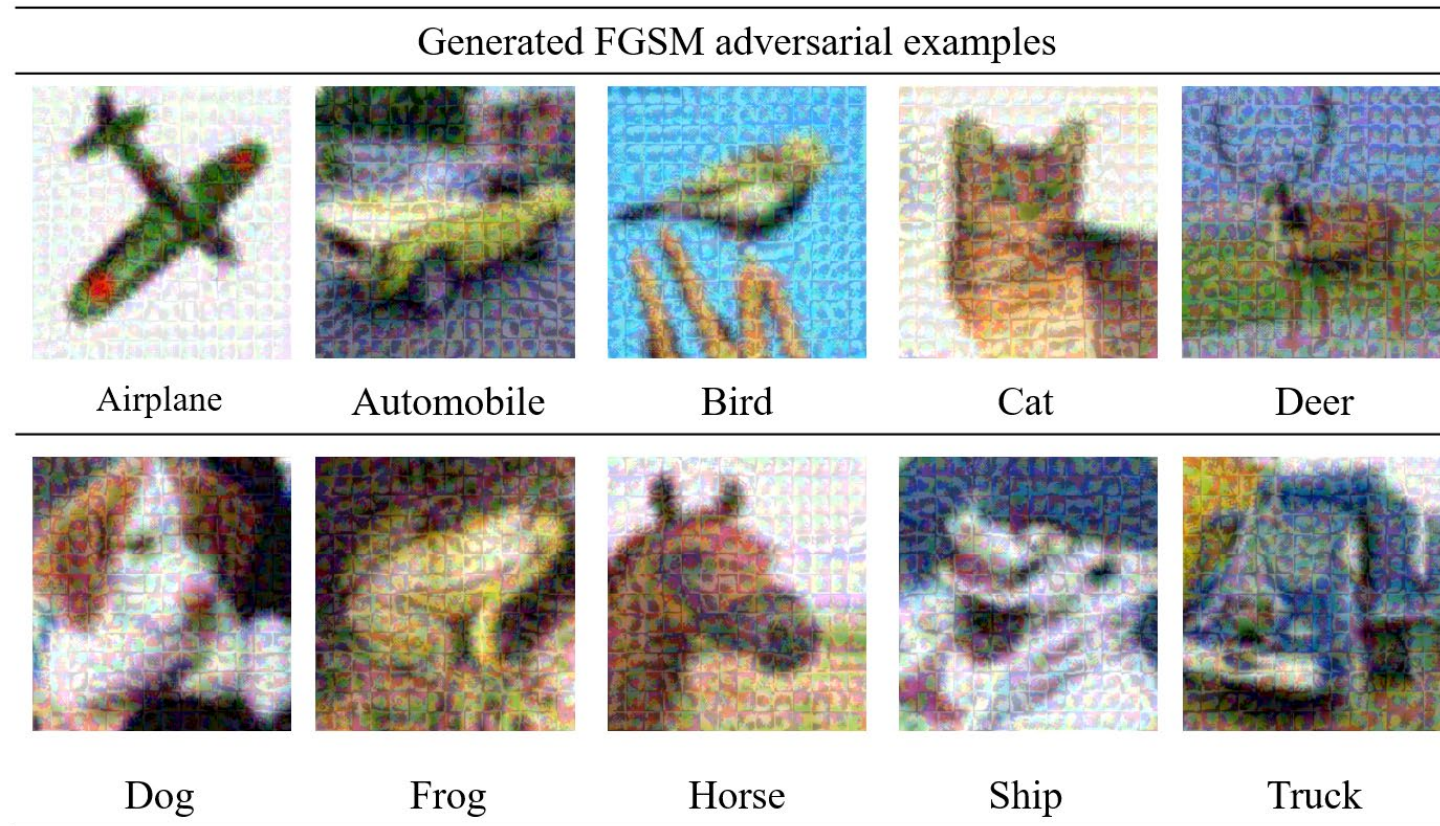
## DeiT's security verification method

- 1) **Fine-tune** of DeiT
- 2) **Weight extract** of fine-tuned DeiT
- 3) **Generate of FGSM adversarial examples** through DeiT weights
- 4) **Insert of generated adversarial examples** into fine-tuned DeiT and **performance analysis**



# • Experiment & Evaluation

## ▪ Generating adversarial examples through FGSM



<Fig7> Generated FGSM adversarial examples

- Generated by FGSM adversarial examples are composed of 100 images per class (A total of 1,000 images)

# • Experiment & Evaluation

- **DeiT performance verification (Accuracy)**

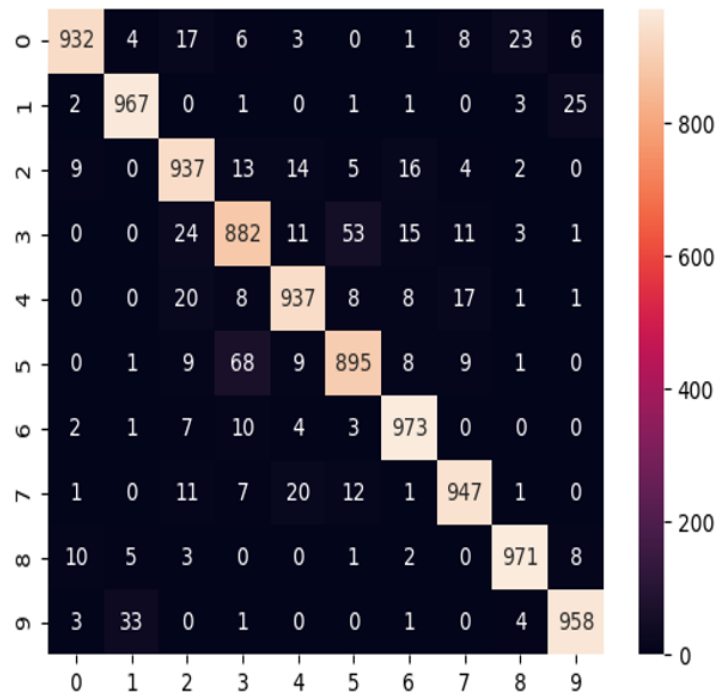
Model	Attack	Dataset	Precision	Recall	F1-Score	Accuracy
DeiT	X	Cifar-10	0.9399	0.9399	0.9399	<b>0.9399</b>
	O	Adversarial examples	0.1050	0.1050	0.1049	<b>0.1050</b>

<Table 1> DeiT Verification Results

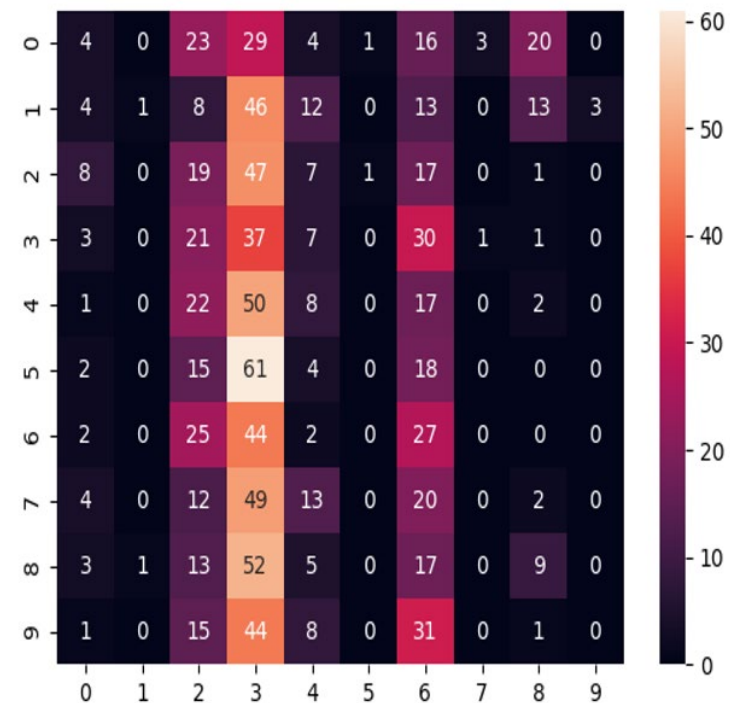
- **83.49% drop** compared to the normal dataset in case of adversarial attack
- We verified that **adversarial attack on DeiT is fatal**

# • Experiment & Evaluation

## ▪ DeiT performance verification (Confusion Matrices)



<Fig8> Original Dataset's Verification (DeiT)



<Fig9> Adversarial Examples Verification (DeiT)

- Adversarial example confusion matrix in Figure 9 : **Proper classification was not achieved.**

# • Conclusion & Future Work

---

- Currently, performance-oriented research is actively in progress in the field of AI.
- However, as new deep learning models develop, research from a security perspective must also be conducted.
- Through this experiment, we suggest that DeiT's security problem exists.

**(DeiT accuracy drop of 83.49% in FGSM Attack)**

- **We suggest**

*1) DeiT's defensive limitations exist*

*2) Need to address DeiT's security vulnerabilities*

- **Future Work**

*1) Create a robust DeiT Model (against adversarial attack)*

---

# Thank you

Inpyo Hong

[hip9863@gachon.ac.kr](mailto:hip9863@gachon.ac.kr)

Dept. of Computer Engineering, Gachon University