Accelerated
Intelligent
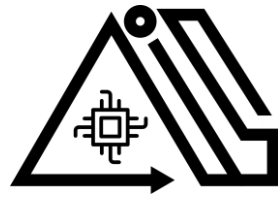Systems Lab.

CVPR
JUNE 17-21, 2024
SEATTLE, WA

# PeerAiD:
## Improving Adversarial Distillation from a Specialized Peer Tutor

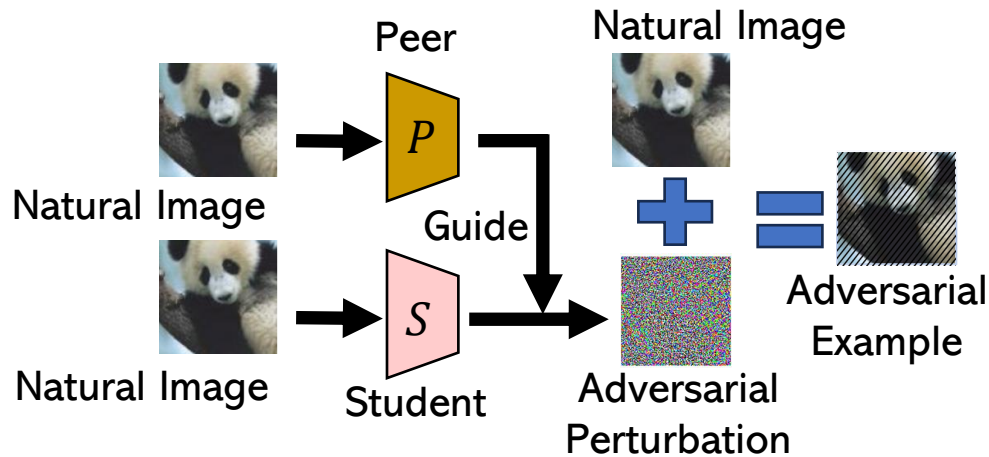**Jaewon Jung**, Hongsun Jang, Jaeyong Song, and Jinho Lee

Department of Electrical and Computer Engineering, Seoul National University
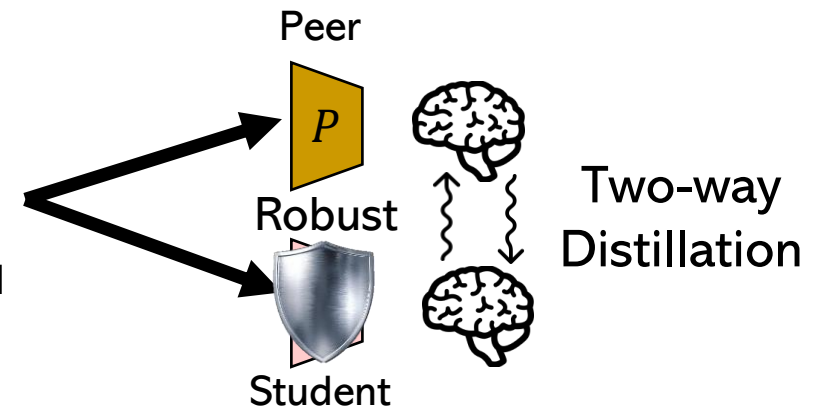
# Overview

- PeerAiD proposes using the peer, which interactively learns with the student during adversarial distillation.
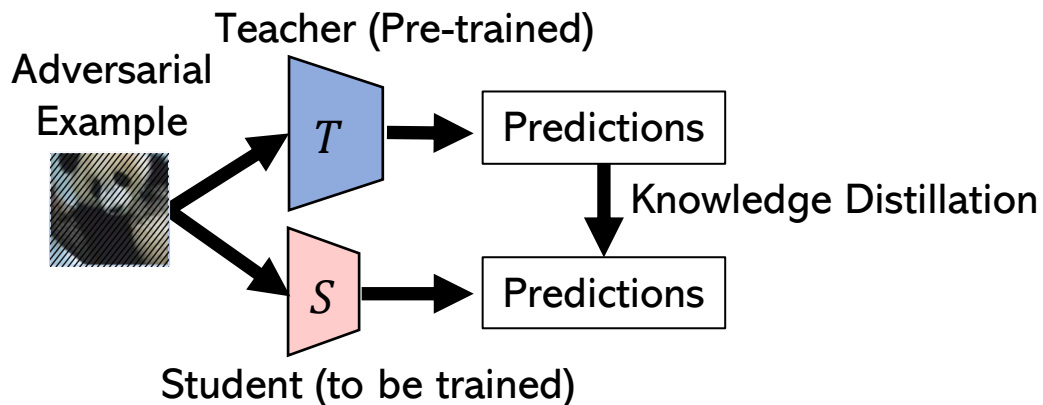
① Adversarial example generation

② Weight optimization



- Conventional AD



- PeerAiD

# Overview

- **PeerAiD proposes using the peer, which interactively learns with the student during adversarial distillation.**
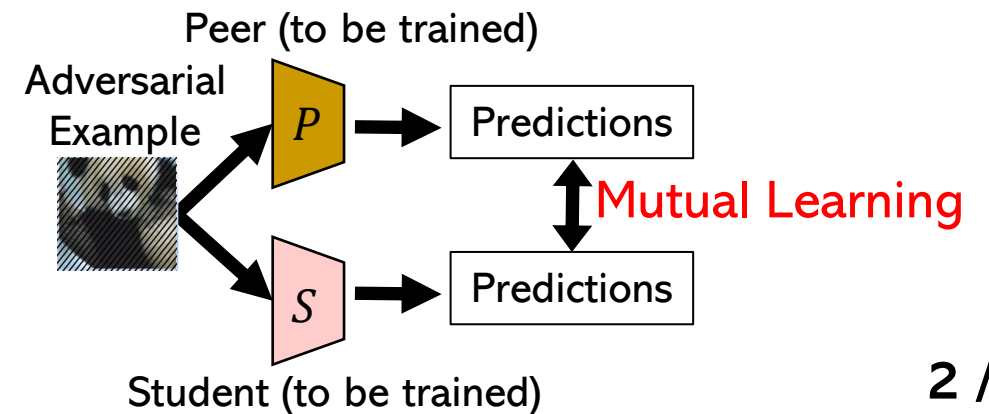
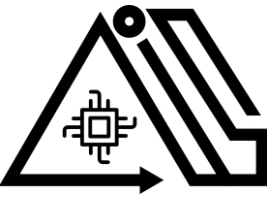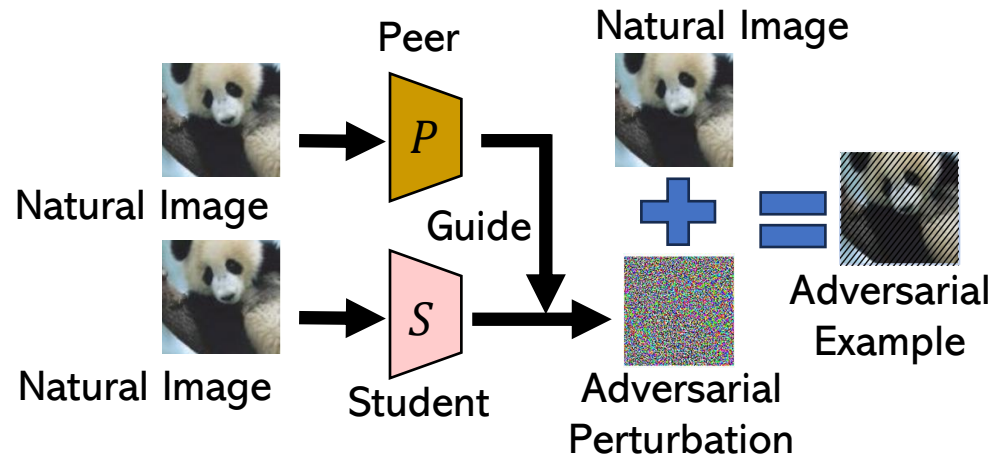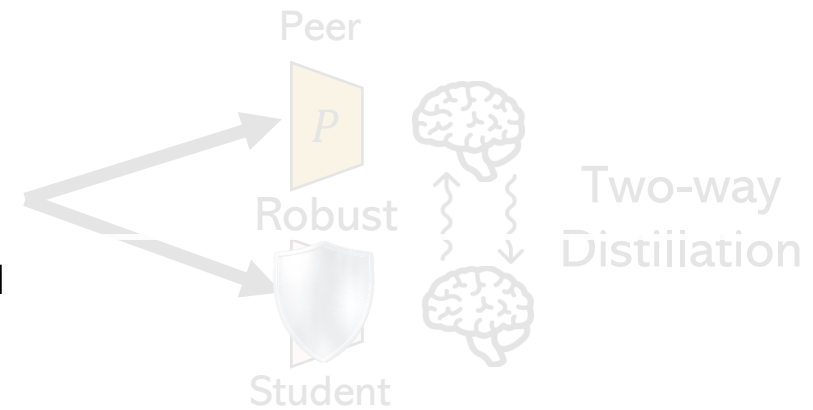① Adversarial example generation



② Weight optimization

Peer

Robust

Two-way Distillation

Student

- Conventional AD

Teacher (Pre-trained)

Adversarial Example

T → Predictions

Knowledge Distillation

S → Predictions

Student (to be trained)

- PeerAiD

Peer (to be trained)

Adversarial Example

P → Predictions

Mutual Learning

S → Predictions

Student (to be trained)

# Overview

- **PeerAiD proposes using the peer, which interactively learns with the student during adversarial distillation.**

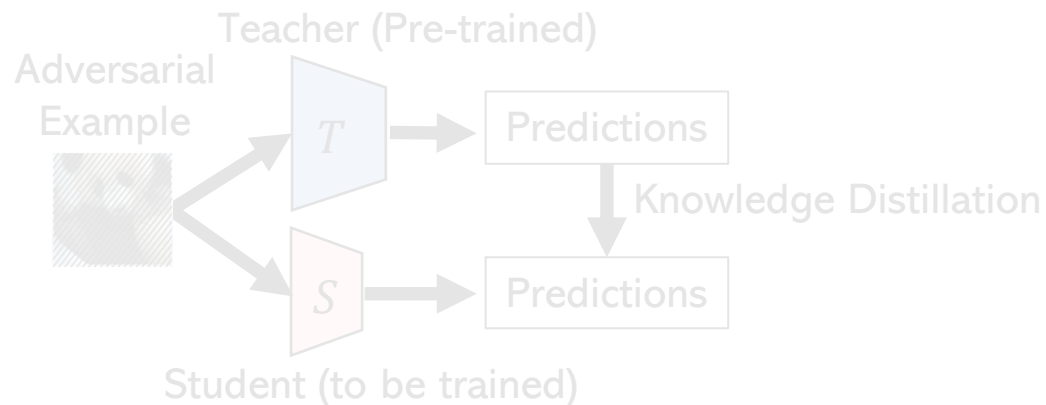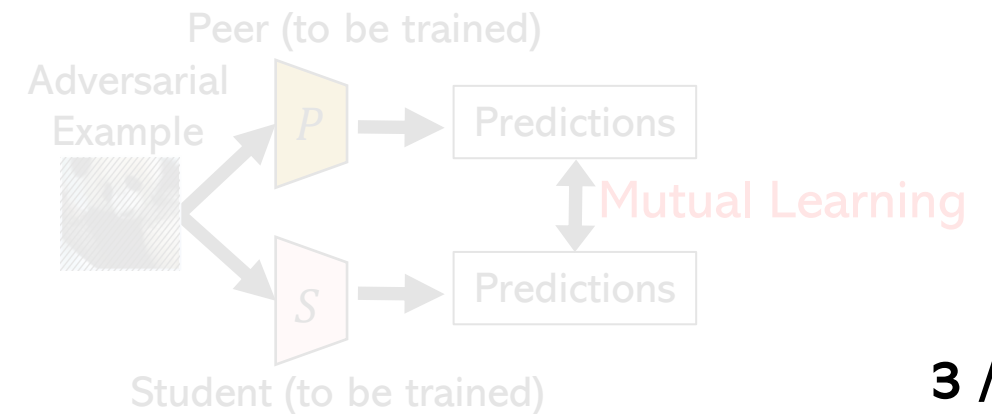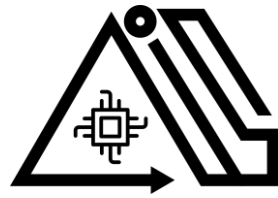① Adversarial example generation

② Weight optimization



Peer

Natural Image

P

Natural Image

Guide

S

Natural Image

Student

Natural Image

Adversarial Perturbation

Adversarial Example

Peer

P

Robust

Student

Two-way Distillation

- Conventional AD

Adversarial Example

Teacher (Pre-trained)

T

Predictions

Knowledge Distillation

S

Predictions

Student (to be trained)

- PeerAiD

Adversarial Example

Peer (to be trained)

P

Predictions

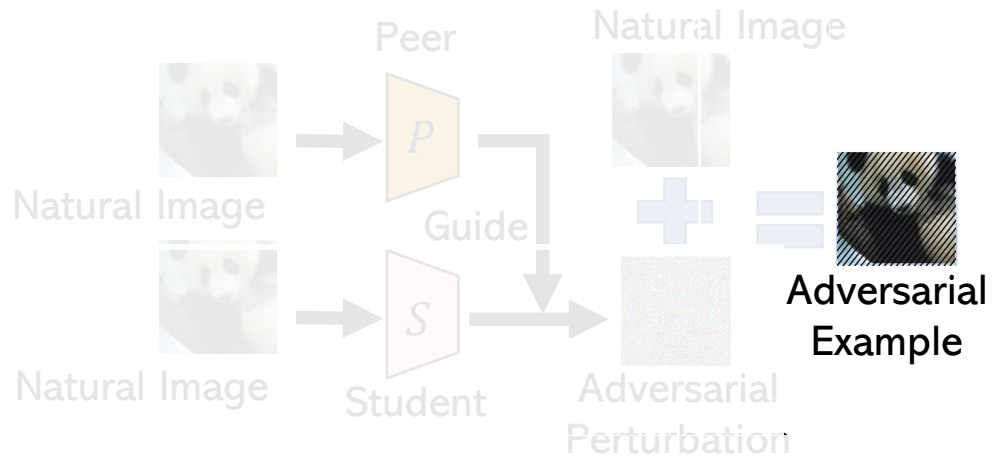Mutual Learning

S

Predictions

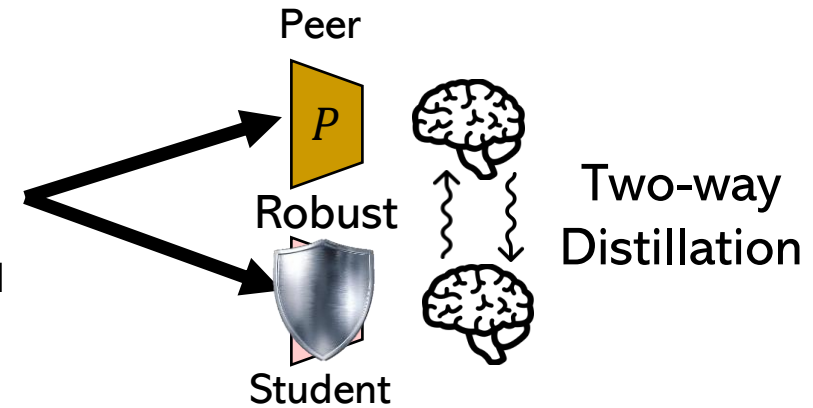Student (to be trained)

# Overview

- PeerAiD proposes using the peer, which interactively learns with the student during adversarial distillation.

① Adversarial example generation

② Weight optimization

Peer

Natural Image

Natural Image

Peer

Guide

Robust

Two-way
Distillation

$P$

$S$

Student

Adversarial
Example

Adversarial
Perturbation

Student

- Conventional AD

Adversarial
Example

Teacher (Pre-trained)

$T$ → Predictions

Knowledge Distillation

$S$ → Predictions

Student (to be trained)

- PeerAiD

Peer (to be trained)

Adversarial
Example

$P$ → Predictions

Mutual Learning

$S$ → Predictions

Student (to be trained)

# Overview

- PeerAiD proposes using the peer, which interactively learns with the student during adversarial distillation.

# Overview

- PeerAiD proposes using the peer, which interactively learns with the student during adversarial distillation.
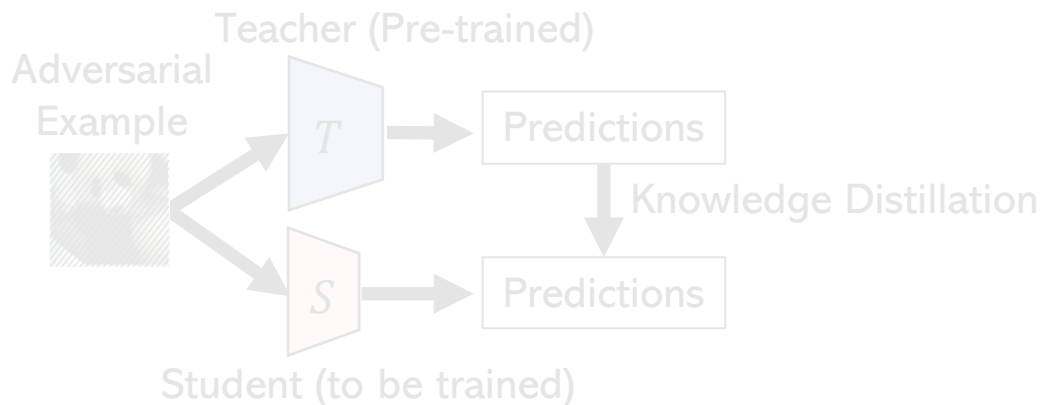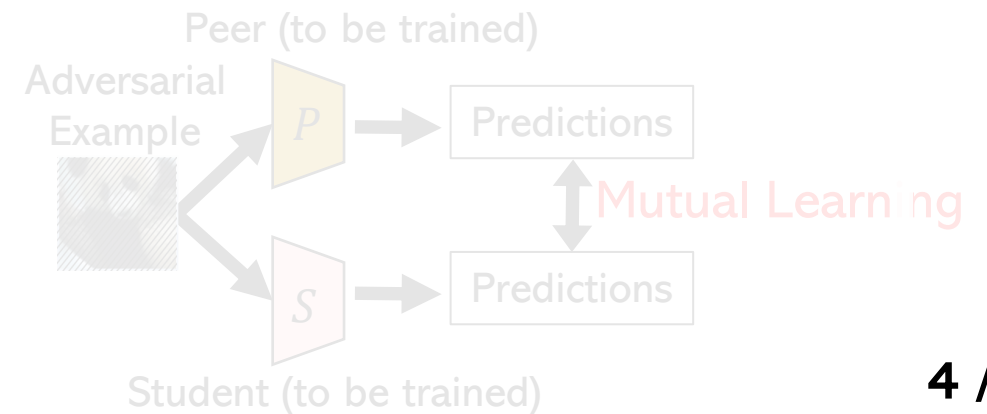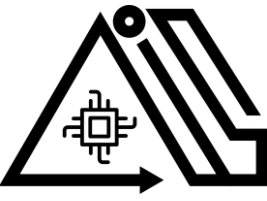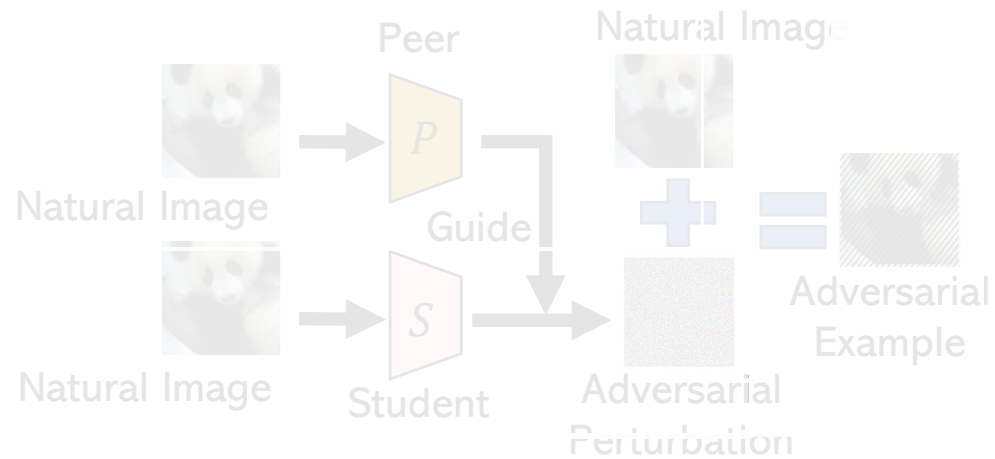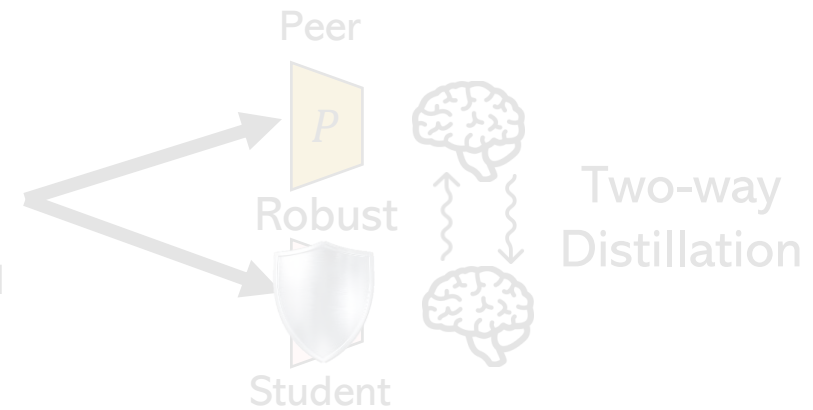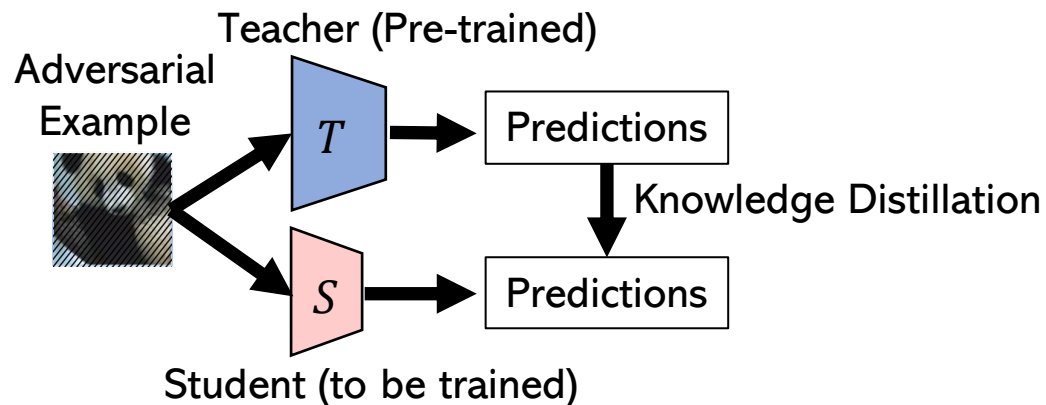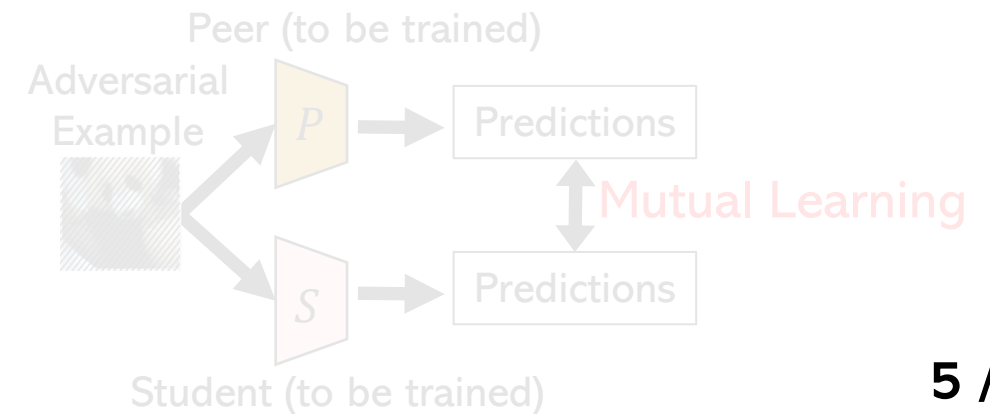
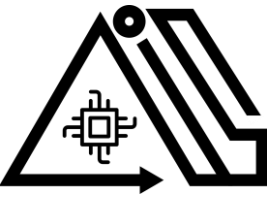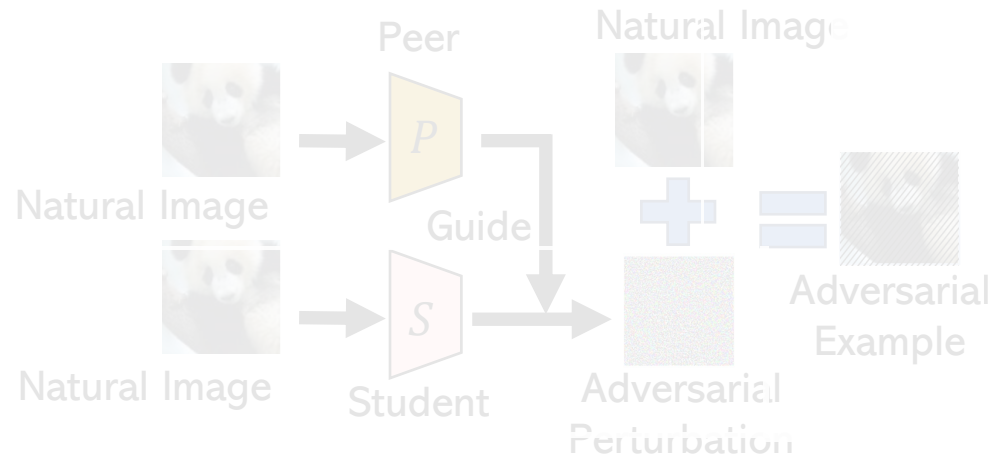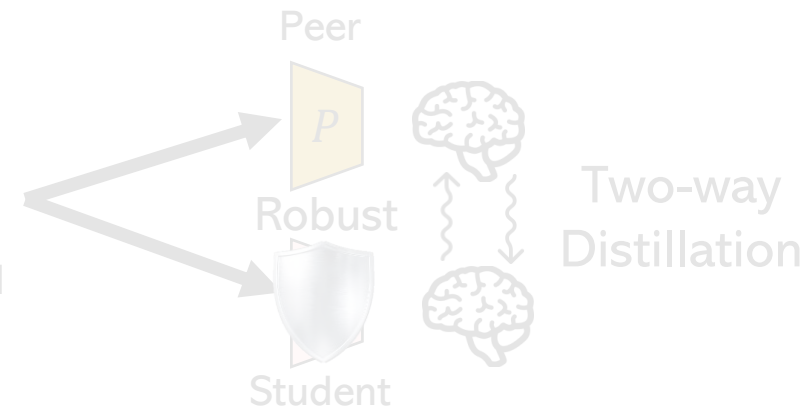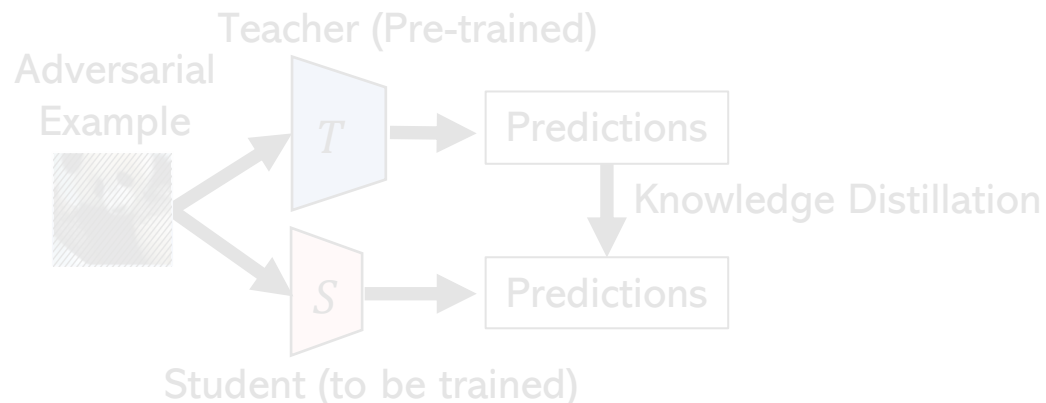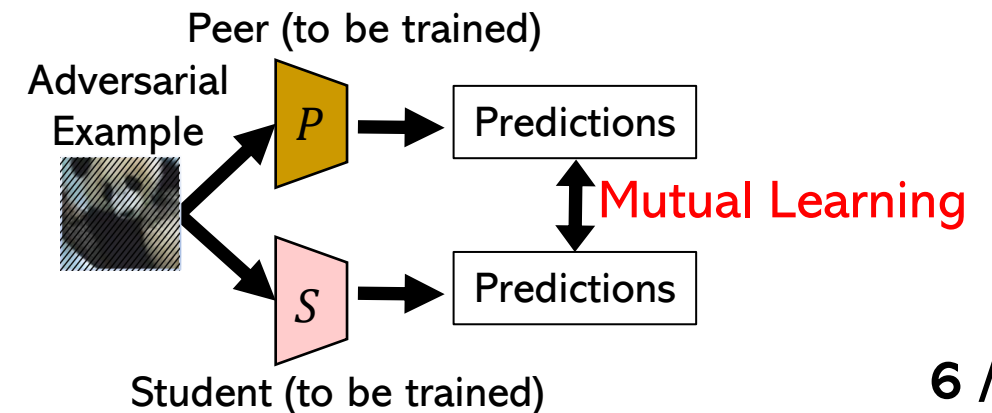① Adversarial example generation
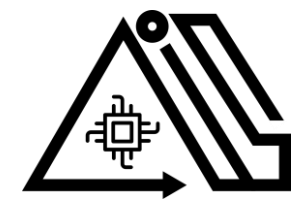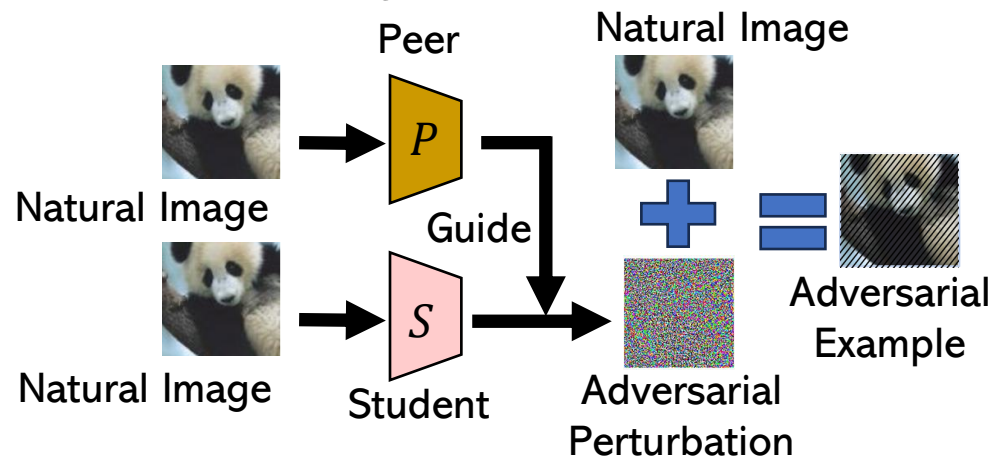
② Weight optimization



- Conventional AD



- PeerAiD

# Background

- **Limitations in the previous works**
  - The pretrained robust teacher model keeps losing its ability to defend against adversarial examples ($x_S^*$) of the student model.



Teacher (Pre-trained)

Student (to be trained)

$T$

$S$

Adversarial Example

Mobile Phone / Embedded Device

Knowledge Distillation

Predictions

Predictions

Attack

Attack

Adversarial Example ($x_T^*$)   Adversarial Example ($x_S^*$)

At later epochs, $x_S^*$ fools the teacher.

Robust Acc. (%) against $x_S^*$ of the pretrained teacher T

Robust Acc. (%) against $x_T^*$ of the pretrained teacher T

Degradation

Test Robust Acc. (%)

Epochs

[ CIFAR-100 result ]

- **Peer Tutoring**
  - Adversarial Example Generation.
    - The student model uses the predictions of the peer model as guidance.



Natural Image ($x$)

Peer Model

: Forward    : Input    : Backward

$p_P(y|x)$

$KL$

Student Model

Adversarial Image ($x_S^*$)

$p_S(y|x)$

Maximize

Natural Image    Adversarial Perturbation

# Proposed Method

- **Peer Tutoring**
  - ## Weight Optimization.
    - The student and the peer transfer their own knowledge to each other.



Natural Image ($x$)

Peer Model

: Forward    : Input    : Backward

$p_P(y|x_S^*)$

$p_P(y|x_S^*)$

$KL$

$y$

Student Model

Adversarial Image ($x_S^*$)

$p_S(y|x_S^*)$

$p_S(y|x_S^*)$    $p_S(y|x)$

$KL$

Minimize

$p_S(y|x_S^*)$

$CE$

$CE$

# Proposed Method

- **Peer Tutoring**
  - Weight Optimization.
    - The student and the peer transfer their own knowledge to each other.

# Proposed Method

- **Peer Tutoring**
  - Weight Optimization.
    - The student and the peer transfer their own knowledge to each other.



Natural Image ($x$)

Peer Model

Adversarial Image ($x_S^*$)

Student Model

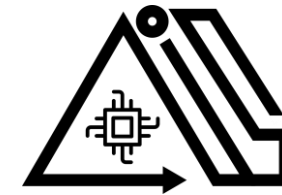: Forward  : Input  : Backward

$p_P(y|x_S^*)$

$p_P(y|x_S^*)$

$KL$

$p_S(y|x_S^*)$

$p_S(y|x_S^*)$

$p_S(y|x)$

$KL$

$p_S(y|x_S^*)$

Minimize

Consistent guidance
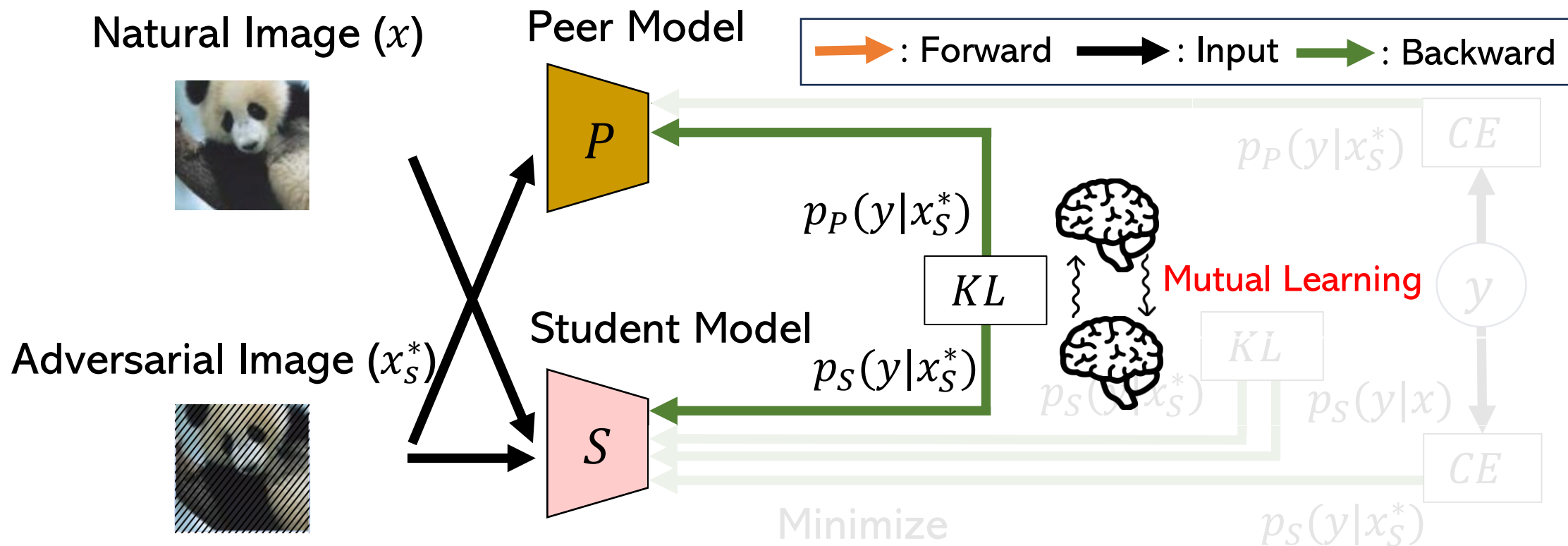for both the peer and the student.

$CE$

$y$

$CE$

$P$

$S$

# Proposed Method

- **Peer Tutoring**
  - ## Weight Optimization.
    - The student and the peer transfer their own knowledge to each other.



Natural Image ($x$)

Peer Model

: Forward    : Input    : Backward

$p_P(y|x_S^*)$    $CE$

$P$

$p_P(y|x_S^*)$

$KL$

Student Model

Adversarial Image ($x_S^*$)

$S$

$y$

The peer learns $x_S^*$ directly.
⇒ It becomes a specialist defender.

# Proposed Method

- **Peer Tutoring**
  - Weight Optimization.
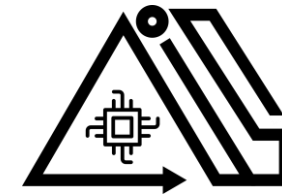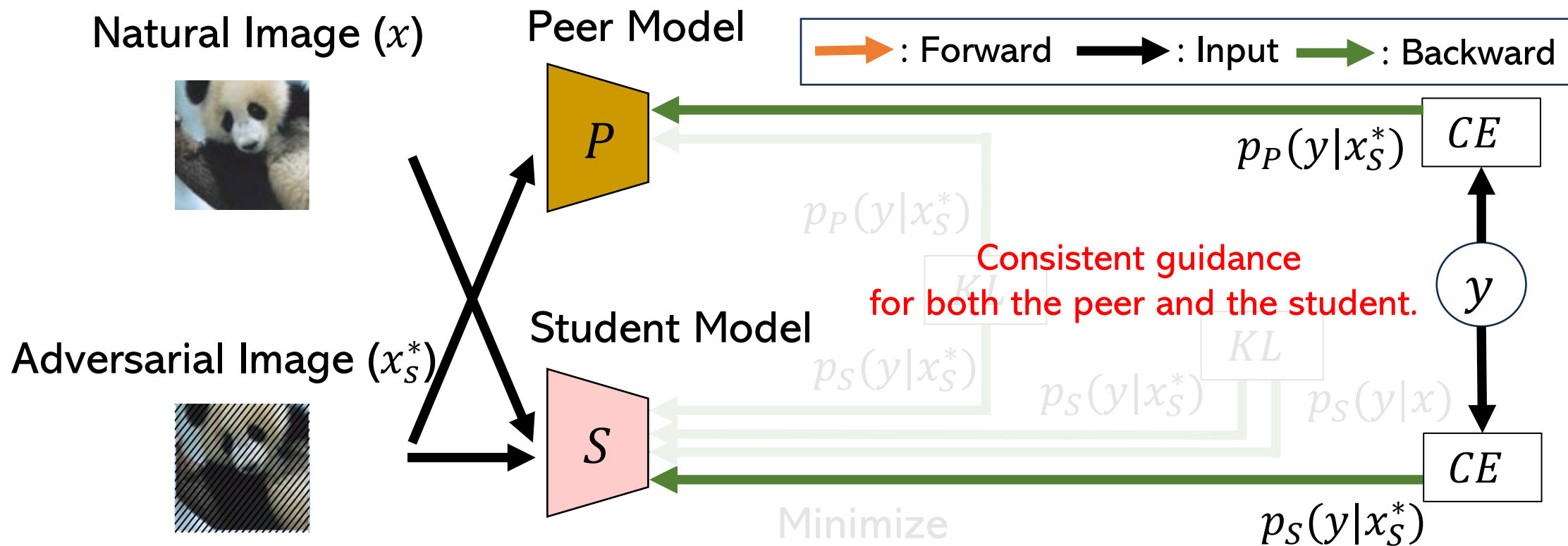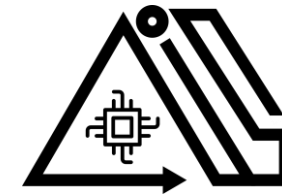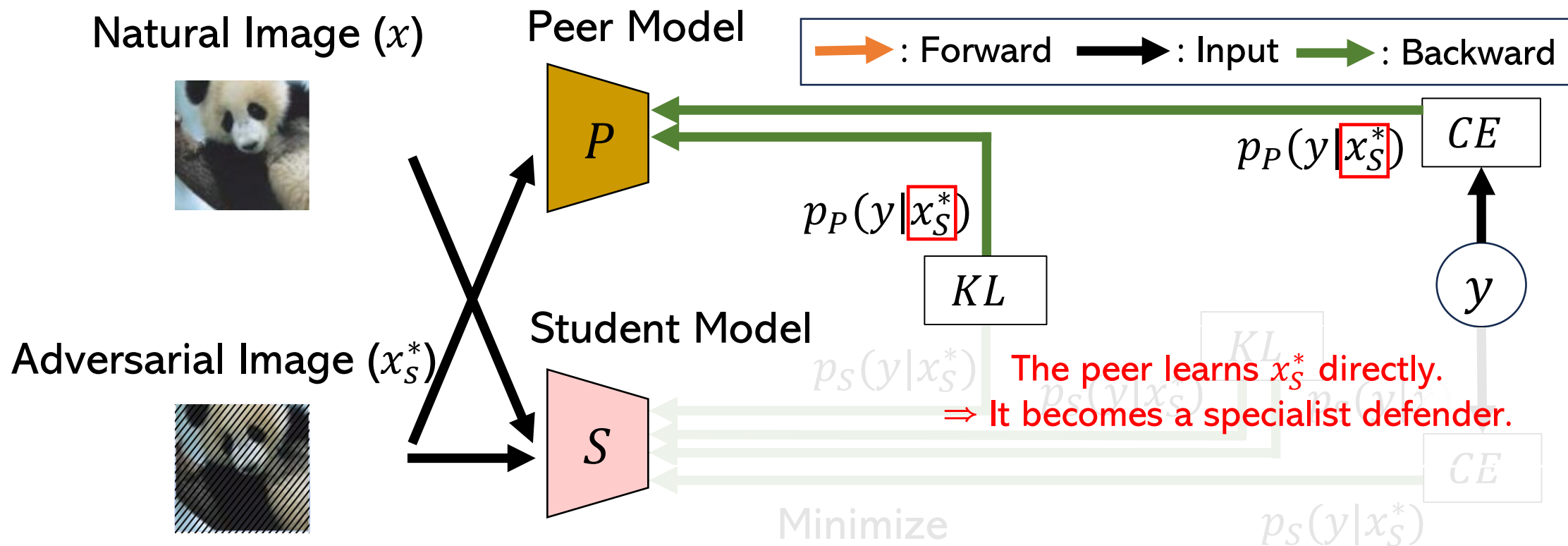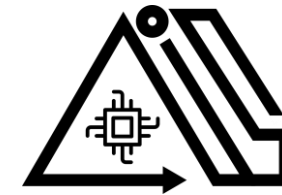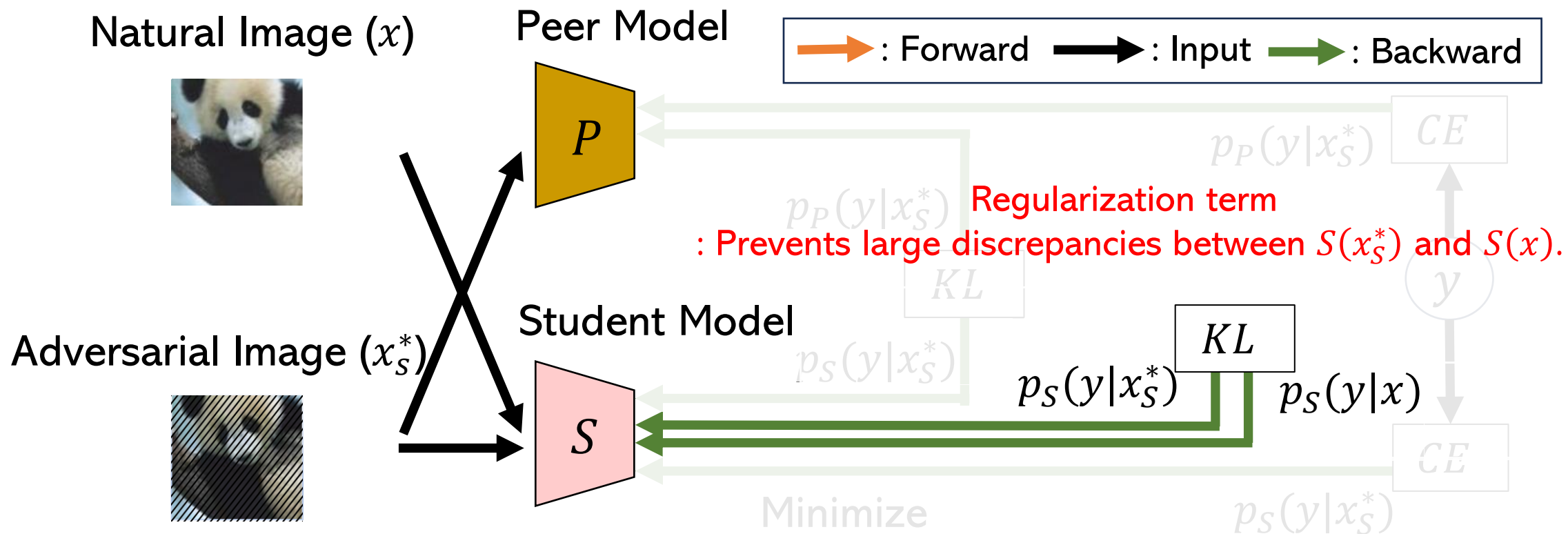    - The student and the peer transfer their own knowledge to each other.



Natural Image ($x$)

Peer Model

: Forward    : Input    : Backward

$p_P(y|x_S^*)$    $CE$

$P$

$p_P(y|x_S^*)$    Regularization term
: Prevents large discrepancies between $S(x_S^*)$ and $S(x)$.

$KL$

Adversarial Image ($x_S^*$)

Student Model

$p_S(y|x_S^*)$

$S$

$KL$

$p_S(y|x_S^*)$    $p_S(y|x)$

$y$

$CE$

Minimize    $p_S(y|x_S^*)$

# Experimental Results

- ## TinyImageNet result
  - PeerAiD shows the highest AutoAttack robust accuracy compared to other baselines, while also providing higher clean accuracy.

# Experimental Results

- ## **Characteristic of the peer model**
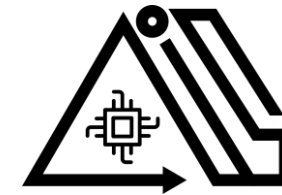
  ① Specialist who defends against adversarial examples of the student model.

  - No tradeoff between the robustness and clean accuracy.
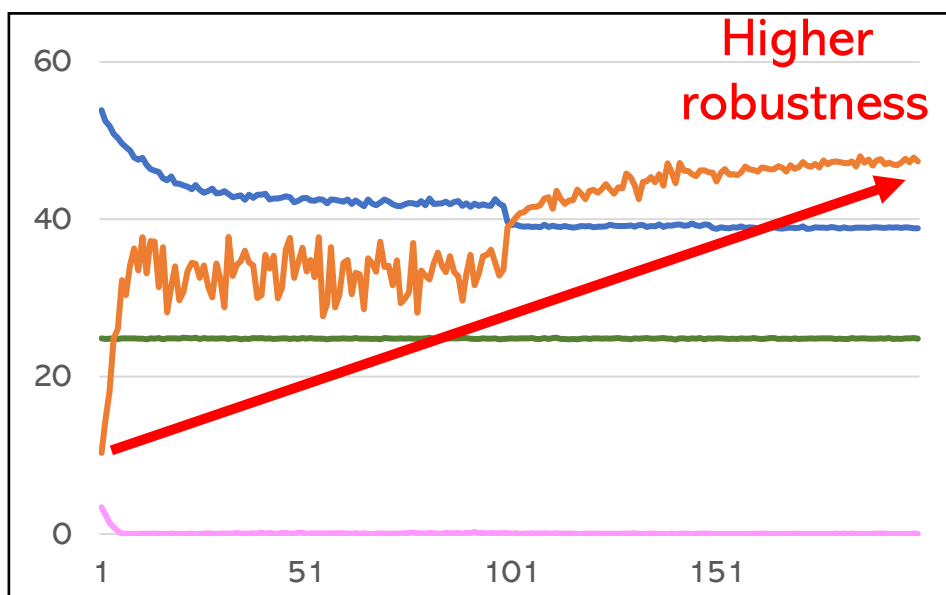
  ② High clean accuracy

— Robust Acc. (%) against $x_S^*$ of the peer model P
— Robust Acc. (%) against $x_P^*$ of the peer model P
— Robust Acc. (%) against $x_S^*$ of the pretrained teacher T
— Robust Acc. (%) against $x_T^*$ of the pretrained teacher T



Higher robustness

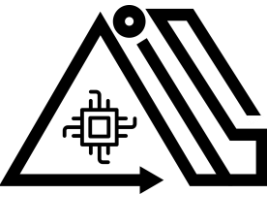| $f$ | $Rob_f(x_f^*)$ | $Rob_f(x_s^*)$ | | Student's Robust acc. |
|---|---|---|---|---|
| Peer tutor | 0.00 | <u>69.19</u> | ⇒ | <u>**29.69**</u> |
| Pretrained robust teacher | 24.15 | 39.46 | ⇒ | 24.48 |

[ CIFAR-100, ResNet-18 result ]

- $Rob_f(\cdot)$ : the robust accuracy of $f$.
- $S$ : the student model.
- Peer's Clean acc : 75.63 > 75.48 (Naturally trained)

# Conclusion

- We propose a novel online adversarial distillation method, PeerAiD

- The peer model specializes in defending against the student model's attack samples.

- PeerAiD improves AA robust accuracy by 1.66%p and clean accuracy by 4.72%p.

# Thanks!

Jaewon Jung @ Seoul National University

Email: jungjaewon@snu.ac.kr

Code: https://github.com/jaewonalive/PeerAiD