

## A. Supplementary Experiments for Section 3.1

This section supplements Section 3.1 of the main paper by providing additional experimental details and results. We present the complete experimental setup described in Section 3.1, results for perturbation size of 8/255 in Table 5, and comprehensive results across five seeds for varying perturbation sizes as illustrated in Figures 6 and 7.

To explore the phenomenon of the transferability of adversarial perturbations following catastrophic overfitting, we engaged in an experimental study based on the settings established by He et al. [11]. We initiated our study by training a ResNet18 [10] model on the CIFAR-10 dataset using the Fast Gradient Sign Method Adversarial Training (FGSM-AT) [7] for 100 epochs. The training was conducted with perturbation sizes ( $\epsilon$ ) set to 8/255 and 16/255, employing a Stochastic Gradient Descent (SGD) optimizer with an initial learning rate of 0.1. The learning rate was programmed to diminish by a factor of 0.1 upon reaching the 80th and 90th epochs. Additionally, the model training incorporated a batch size of 128, and the images underwent preprocessing, which included padding of 4 pixels on each side, followed by random cropping and horizontal flipping. To precisely replicate the conditions leading to CO, we adopted zero initialization for generating adversarial samples and set the weight decay to zero. This setup was chosen to maintain consistency with He et al. [11] and to ensure the stable reproduction of CO, thereby facilitating a clear examination of the transferability of adversarial perturbations under these conditions.

We delve into the specifics of the experimental outcomes for each seed, as illustrated in Figures 6 and 7, to shed light on the underlying dynamics of  $P_{abnormal}$  and PGD accuracy in relation to catastrophic overfitting. For all five seeds, we observed a gradual increase in  $P_{abnormal}$  during the initial stages of training. However, a striking observation was made at the point of CO, where PGD accuracy plummeted to approximately 0, underscoring a sudden and severe degradation in the model’s ability to counter adversarial attacks. Correspondingly,  $P_{abnormal}$  experienced a sharp escalation, reinforcing the strong linkage between the onset of CO and the dramatic increase in  $P_{abnormal}$ . This pattern was consistent across different seeds. The detailed analysis for each seed further corroborates the significant impact of CO on the transferability of adversarial perturbations.

## B. Pseudocode of the LIET Algorithm for Section 3.3

This section supplements Section 3.3 of the main paper by presenting the pseudocode for the LIET algorithm (Algorithm 1). To enhance clarity, the pseudocode simplifies certain operations. For instance, an implementation detail con-

---

### Algorithm 1: LIET: Label Information Elimination Training

---

**Input:** A classifier  $f_\theta$  with loss function  $\mathcal{L}$ ; Dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ ; Perturbation magnitude  $\epsilon$ ; Gray image  $x_{gray}$ ; Hyperparameter  $\lambda$ ; Number of epochs  $E$ .

**Output:** Robust model parameters  $\theta$

```

1 for  $e = 1$  to  $E$  do
2   for each class  $c$  in dataset do
3      $LI_c = \epsilon \cdot \text{sign}(\nabla_{x_{gray}} \mathcal{L}(f(x_{gray}; \theta), c))$ 
      {Generate class-specific label information}
4   end
5   for each batch  $\mathcal{B} = (x, y) \in \mathcal{D}$  do
6     if  $\text{random}() < 0.5$  then
7        $x' = x + LI_y$  {Randomly add label
        information}
8     else
9        $x' = x - LI_y$  {Randomly subtract
        label information}
10    end
11     $\delta_x = \epsilon \cdot \text{sign}(\nabla_{x'} \mathcal{L}(f(x'; \theta), y))$  {Generate
      adversarial perturbation}
12     $x_{adv} = x + \delta_x$  {Create adversarial
      example}
13     $\text{loss}_1 = \mathcal{L}(f(x_{adv}; \theta), y)$  {Standard
      adversarial training loss}
14     $\text{loss}_2 = \lambda \cdot \mathcal{L}_{JSD}(f(x; \theta), f(x_{adv}; \theta))$  {JS
      divergence for smoother loss surface}
15     $\text{total\_loss} = \text{loss}_1 + \text{loss}_2$  {Combined loss
      function}
16  end
17   $\theta = \theta - \eta \cdot \nabla_{\theta} \text{total\_loss}$  {Update model
    parameters}
18 end

```

---

cerns the clipping of adversarial perturbations. Specifically, for a perturbation budget of  $\epsilon = 8/255$ , we clip the perturbation to stay within this bound. For larger budgets, however, we adopt the strategy from [3] and omit the clipping step to generate stronger adversaries. For a comprehensive implementation, we refer the reader to the provided source code.

## C. Experiment Details for Section 4.1

This section provides supplementary information to Section 4.1 of the main paper. Here, we present detailed experimental configurations and parameters that were utilized in our study but were omitted from the main paper for brevity.

Our research conducted experiments on three widely recognized datasets to evaluate the robustness against adversarial attacks, namely CIFAR-10, CIFAR-100, and

Input $x$	Perturbation: 16/255		Perturbation: 8/255	
	$P_{abnormal}$ (%)	$P_{dominate}$ (%)	$P_{abnormal}$ (%)	$P_{dominate}$ (%)
Uniform Gray (0)	$42.20 \pm 32.55$	$20.29 \pm 28.62$	$49.89 \pm 1.34$	$10.08 \pm 0.16$
Uniform Gray (0.5)	$81.81 \pm 20.00$	$44.47 \pm 27.44$	$53.52 \pm 4.13$	$10.52 \pm 0.62$
Training Mean	$81.36 \pm 20.47$	$45.45 \pm 28.42$	$53.84 \pm 4.31$	$10.55 \pm 0.68$
Uniform Noise	$50.70 \pm 2.35$	$11.08 \pm 1.72$	$48.67 \pm 0.93$	$9.87 \pm 0.10$
Test Sample 0	$73.08 \pm 18.95$	$26.06 \pm 15.03$	$51.18 \pm 0.49$	$10.12 \pm 0.13$

Table 5. Transferability of label information for different inputs on the CIFAR-10 dataset with perturbation sizes of 16/255 and 8/255.

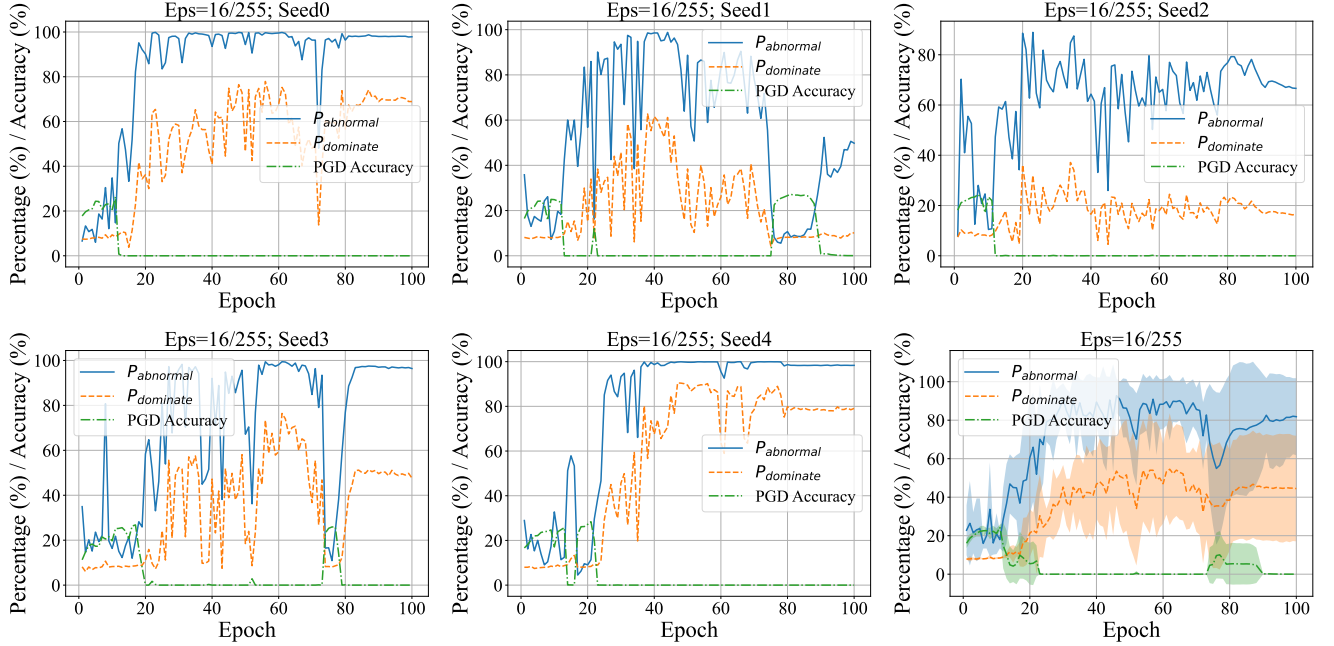


Figure 6. Transferability of label information for uniform gray image (value 0.5) on the CIFAR10 dataset with a perturbation size of 16/255.

Tiny ImageNet. For CIFAR-10 and CIFAR-100 datasets, we employed both the ResNet-18 architecture [10] and WideResNet-28-10 [29] as the network architectures, while for Tiny ImageNet, we opted for PreActResNet18 due to its enhanced performance on more complex datasets.

For the CIFAR-10, CIFAR-100, and Tiny ImageNet datasets, we carved out validation sets comprising 1000, 1000, and 2000 images, respectively, from the training data. During the training phase, we evaluated the model’s performance on these validation sets using the PGD-10 accuracy metric. The model that achieved the highest accuracy on the validation set was selected as the final model. This validation strategy was consistently applied across all compared algorithms to maintain uniformity in model evaluation.

We set a batch size of 128 and applied a series of preprocessing steps on the images. These steps included padding the images with 4 pixels on each side, followed by random

cropping and horizontal flipping to augment the dataset and improve model generalization.

We utilized the Stochastic Gradient Descent (SGD) as our optimization algorithm, with an initial learning rate set at 0.1, a weight decay parameter of  $5e-4$ , and momentum of 0.9. The training process was conducted over 100 epochs, incorporating a OneCycleLR scheduler to adjust the learning rate dynamically. To stabilize the training process, we implemented a Weight Averaging (WA) [12] technique with a  $\tau$  value of 0.9995. Each experiment was replicated three times under different random seeds to ensure the reliability of our results. For perturbation magnitude of 8/255, we set  $\lambda$  values at 100, 200, and 100, respectively. Furthermore, we employed non-uniform label smoothing values of 0.6 for both CIFAR-10 and CIFAR-100, and 0.8 for Tiny ImageNet, to fine-tune the model’s performance across diverse datasets. For perturbation magnitude of 16/255, we set  $\lambda$  to

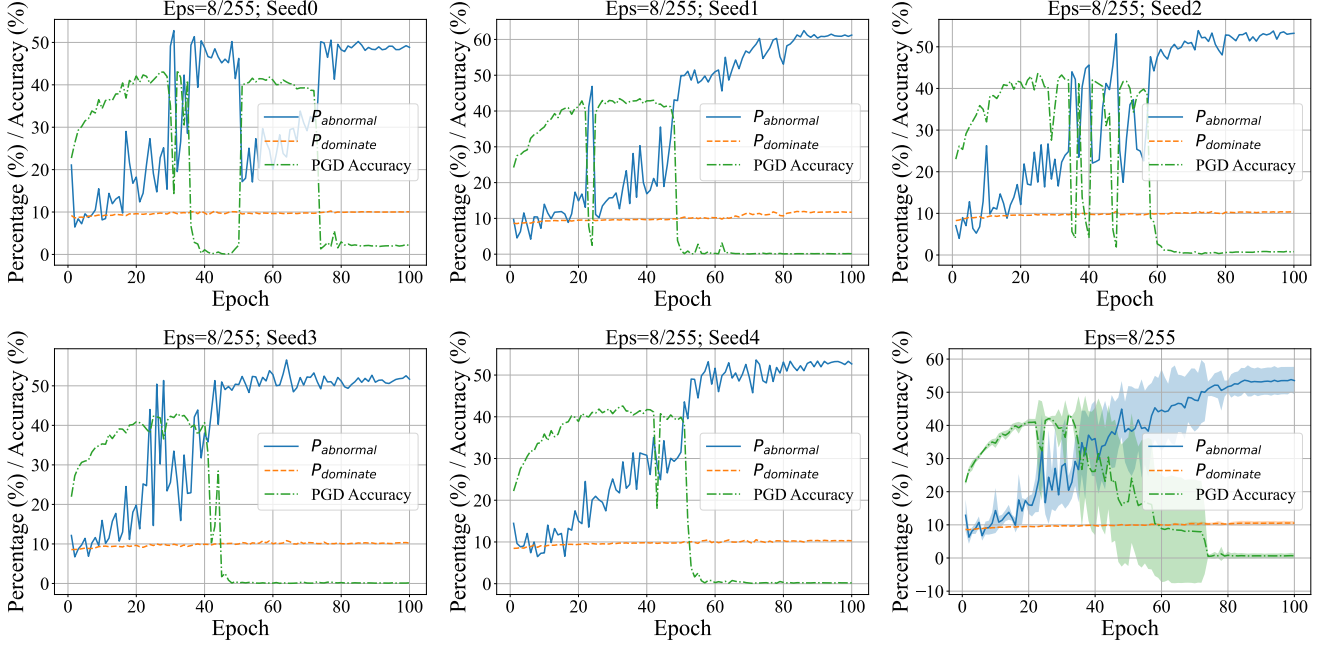


Figure 7. Transferability of label information for uniform gray image (value 0.5) on the CIFAR10 dataset with a perturbation size of 8/255.

20 and label smoothing value to 0.4.

As highlighted in our paper, initializing the training samples,  $\mathbf{x}$ , by either adding or subtracting  $\mathbf{LI}_c$  proved effective in diminishing the label information,  $y$ , from the generated perturbation,  $\delta_{\mathbf{x}}$ . This initialization strategy was employed randomly to enhance the unpredictability of our defense mechanism against adversarial inputs.

In line with our strategy to boost diversity within the model, we randomly substituted 10% to 50% of the elements in  $\mathbf{LI}_c$  with values uniformly distributed between  $-\epsilon$  and  $\epsilon$  when the perturbation size was 8/255. For a perturbation size of 16/255, we randomly substituted 0% to 100% of the elements in  $\mathbf{LI}_c$  with values uniformly distributed between  $-2\epsilon$  and  $2\epsilon$ . This approach was aimed at enriching the robustness of our model against adversarial attacks. To ensure our model adapts to evolving adversarial tactics, we updated  $\mathbf{LI}_c$  at different intervals depending on the dataset: every 10 batches for CIFAR10 and every 20 batches for CIFAR100 and Tiny-ImageNet, aligning with our strategy to maintain model resilience over time.

To evaluate the robustness of our model, we subjected it to several adversarial attack methods, including PGD [16] and AutoAttack (AA) [2]. We varied the number of iterations for PGD attacks to 10, 20, and 50, which are henceforth referred to as PGD-20, and PGD-50, respectively.

The maximum allowed perturbation was set to 8/255 and 16/255 for CIFAR-10 and CIFAR-100 datasets, respectively. For Tiny-ImageNet, we used perturbation bounds of

8/255 and 12/255. These dataset-specific perturbation settings reflect the different sensitivity levels of each dataset to adversarial attacks and provide a more comprehensive evaluation of our defense mechanism.

## D. Experiment Results for Section 4.2

This section supplements the results presented in Section 4.2 of the main paper by providing more detailed experimental findings. Specifically, Tables 6, 7, and 8 display the average results from three sets of experiments: CIFAR-10 and CIFAR-100 using ResNet-18, and Tiny ImageNet using PreActResNet-18. These results include clean accuracy and robust accuracy measured against PGD-10, PGD-20, PGD-50, and AutoAttack. Additionally, Tables 9 and 10 present experimental results for CIFAR-10 and CIFAR-100 using the larger WideResNet-28-10 architecture. Due to the computational demands of this larger model, these experiments were conducted only once.

For the training costs reported in Tables 6, 7, and 8, we measured training time on an NVIDIA V100 GPU. To calculate computational complexity (PFLOPs), we approximated the backward propagation cost as equivalent to the forward propagation cost. The total FLOPs for each method were determined by calculating the number of forward and backward passes required during training.

It is worth noting that while some methods have identical FLOPs calculations in our tables, their training times differ significantly. This discrepancy arises because a sub-

stantial portion of training time is consumed by parameter updates, and some methods require maintaining computational graphs in memory, which introduces additional overhead not captured in FLOPs measurements alone.

Method	Clean (%) $\uparrow$	PGD-10 (%) $\uparrow$	PGD-20 (%) $\uparrow$	PGD-50 (%) $\uparrow$	AA (%) $\uparrow$	Training Cost
PGD-AT [16]	82.32 $\pm$ 0.39	53.76 $\pm$ 0.18	52.83 $\pm$ 0.11	52.60 $\pm$ 0.13	48.68 $\pm$ 0.12	353.70 min
	70.91 $\pm$ 1.37	37.10 $\pm$ 0.25	28.41 $\pm$ 0.09	25.80 $\pm$ 0.36	20.07 $\pm$ 0.05	59.87 PFLOPs
PGD-AT-WA [24]	82.00 $\pm$ 0.38	54.90 $\pm$ 0.15	54.21 $\pm$ 0.15	54.08 $\pm$ 0.12	50.26 $\pm$ 0.14	358.21 min
	71.10 $\pm$ 0.31	41.09 $\pm$ 0.07	33.22 $\pm$ 0.25	31.17 $\pm$ 0.21	25.09 $\pm$ 0.21	59.87 PFLOPs
N-FGSM [3]	78.79 $\pm$ 0.46	53.78 $\pm$ 0.13	53.28 $\pm$ 0.07	53.17 $\pm$ 0.09	47.97 $\pm$ 0.12	74.15 min
	63.09 $\pm$ 2.38	<b>39.29 <math>\pm</math> 0.45</b>	33.00 $\pm$ 1.05	31.65 $\pm$ 1.36	22.50 $\pm$ 1.33	10.89 PFLOPs
FGSM-RS [28]	73.06 $\pm$ 10.15	48.22 $\pm$ 7.64	47.77 $\pm$ 7.46	47.68 $\pm$ 7.46	42.82 $\pm$ 6.70	74.36 min
	66.21 $\pm$ 7.19	15.11 $\pm$ 5.15	9.33 $\pm$ 6.20	5.42 $\pm$ 5.12	0.00 $\pm$ 0.00	10.89 PFLOPs
FGSM-PGI [13]	80.32 $\pm$ 1.09	56.36 $\pm$ 0.19	55.81 $\pm$ 0.13	55.70 $\pm$ 0.10	49.73 $\pm$ 0.17	99.71 min
	88.06 $\pm$ 0.35	19.02 $\pm$ 0.34	12.96 $\pm$ 0.38	8.78 $\pm$ 0.28	0.11 $\pm$ 0.00	10.89 PFLOPs
FGSM-UAP [20]	79.17 $\pm$ 0.27	56.60 $\pm$ 0.03	56.10 $\pm$ 0.04	55.89 $\pm$ 0.02	49.36 $\pm$ 0.11	108.71 min
	88.07 $\pm$ 0.25	15.56 $\pm$ 1.35	9.72 $\pm$ 1.11	6.03 $\pm$ 0.91	0.03 $\pm$ 0.00	16.33 PFLOPs
NuAT [27]	80.78 $\pm$ 0.55	55.43 $\pm$ 0.08	54.79 $\pm$ 0.04	54.64 $\pm$ 0.05	<b>50.04 <math>\pm</math> 0.13</b>	127.81 min
	<b>91.82 <math>\pm</math> 0.12</b>	14.93 $\pm$ 0.61	7.29 $\pm$ 0.63	3.51 $\pm$ 0.37	0.13 $\pm$ 0.02	16.33 PFLOPs
Grad-Align [1]	78.69 $\pm$ 1.13	53.52 $\pm$ 0.16	52.99 $\pm$ 0.20	52.93 $\pm$ 0.20	47.99 $\pm$ 0.15	228.64 min
	46.93 $\pm$ 26.13	27.72 $\pm$ 12.53	23.09 $\pm$ 9.28	21.82 $\pm$ 8.42	15.43 $\pm$ 3.94	16.33 PFLOPs
FGSM-AT [7]	<b>90.98 <math>\pm</math> 0.46</b>	38.66 $\pm$ 1.92	28.91 $\pm$ 1.36	19.47 $\pm$ 1.65	0.00 $\pm$ 0.00	74.28 min
	79.75 $\pm$ 1.64	15.15 $\pm$ 2.49	10.13 $\pm$ 2.89	5.98 $\pm$ 2.41	0.00 $\pm$ 0.00	10.89 PFLOPs
Free-AT [25]	81.99 $\pm$ 0.95	52.23 $\pm$ 0.15	51.60 $\pm$ 0.09	51.42 $\pm$ 0.08	47.43 $\pm$ 0.14	70.53 min
	89.15 $\pm$ 0.54	31.91 $\pm$ 2.17	21.66 $\pm$ 1.98	12.47 $\pm$ 0.50	0.00 $\pm$ 0.00	10.89 PFLOPs
COAT [14]	83.93 $\pm$ 0.34	53.10 $\pm$ 0.22	52.38 $\pm$ 0.17	52.33 $\pm$ 0.09	37.99 $\pm$ 0.31	90.03 min
	84.58 $\pm$ 5.83	36.27 $\pm$ 6.44	24.64 $\pm$ 6.94	18.17 $\pm$ 1.52	4.58 $\pm$ 3.24	13.61 PFLOPs
GAT [26]	85.12 $\pm$ 0.04	55.03 $\pm$ 0.05	54.23 $\pm$ 0.20	54.05 $\pm$ 0.21	49.39 $\pm$ 0.22	126.75 min
	65.02 $\pm$ 38.91	9.48 $\pm$ 0.79	5.98 $\pm$ 2.88	4.44 $\pm$ 3.94	3.14 $\pm$ 4.41	16.33 PFLOPs
LIET (Ours)	80.61 $\pm$ 0.44	<b>56.70 <math>\pm</math> 0.06</b>	<b>56.14 <math>\pm</math> 0.05</b>	<b>56.08 <math>\pm</math> 0.07</b>	50.01 $\pm$ 0.09	101.29 min
	52.72 $\pm$ 3.51	37.17 $\pm$ 2.11	<b>33.55 <math>\pm</math> 1.59</b>	<b>33.10 <math>\pm</math> 1.44</b>	<b>25.22 <math>\pm</math> 0.41</b>	10.93 PFLOPs

Table 6. Comparison of clean accuracy, robust accuracy and training cost (time in minutes and computation in PFLOPs) on CIFAR-10 using ResNet-18. Each method is evaluated with perturbation sizes of 8/255 (first row) and 16/255 (second row). Best results are highlighted in **bold**.

Method	Clean (%) $\uparrow$	PGD-10 (%) $\uparrow$	PGD-20 (%) $\uparrow$	PGD-50 (%) $\uparrow$	AA (%) $\uparrow$	Training Cost
PGD-AT [16]	57.52 $\pm$ 0.95 48.38 $\pm$ 2.04	29.60 $\pm$ 0.23 17.03 $\pm$ 0.18	28.99 $\pm$ 0.21 12.66 $\pm$ 0.09	28.87 $\pm$ 0.27 11.56 $\pm$ 0.16	25.48 $\pm$ 0.11 9.17 $\pm$ 0.12	353.65 min 59.88 PFLOPs
PGD-AT-WA [24]	56.48 $\pm$ 1.34 45.84 $\pm$ 1.58	32.51 $\pm$ 0.31 22.44 $\pm$ 0.10	32.23 $\pm$ 0.26 18.09 $\pm$ 0.16	32.22 $\pm$ 0.25 17.45 $\pm$ 0.27	26.84 $\pm$ 0.28 12.66 $\pm$ 0.15	358.80 min 59.88 PFLOPs
N-FGSM [3]	54.77 $\pm$ 1.68 39.39 $\pm$ 2.11	30.70 $\pm$ 0.27 19.90 $\pm$ 0.30	30.47 $\pm$ 0.26 16.52 $\pm$ 0.04	30.41 $\pm$ 0.29 15.95 $\pm$ 0.18	25.31 $\pm$ 0.28 11.18 $\pm$ 0.21	74.62 min 10.89 PFLOPs
FGSM-RS [28]	37.59 $\pm$ 15.43 3.34 $\pm$ 0.58	20.94 $\pm$ 9.59 1.73 $\pm$ 0.35	20.86 $\pm$ 9.54 1.62 $\pm$ 0.33	20.84 $\pm$ 9.55 1.61 $\pm$ 0.32	16.70 $\pm$ 7.84 1.05 $\pm$ 0.21	74.62 min 10.89 PFLOPs
FGSM-PGI [13]	56.02 $\pm$ 0.21 57.29 $\pm$ 9.61	32.70 $\pm$ 0.14 3.85 $\pm$ 1.87	32.32 $\pm$ 0.10 2.08 $\pm$ 1.23	32.31 $\pm$ 0.11 1.50 $\pm$ 1.33	26.76 $\pm$ 0.07 0.57 $\pm$ 0.70	99.85 min 10.89 PFLOPs
FGSM-UAP [20]	53.54 $\pm$ 0.49 44.39 $\pm$ 28.33	32.14 $\pm$ 0.05 3.18 $\pm$ 0.66	31.83 $\pm$ 0.04 1.89 $\pm$ 0.10	31.81 $\pm$ 0.05 1.35 $\pm$ 0.47	26.29 $\pm$ 0.03 0.55 $\pm$ 0.67	114.03 min 16.33 PFLOPs
NuAT [27]	57.72 $\pm$ 2.01 62.30 $\pm$ 0.08	25.82 $\pm$ 1.28 10.72 $\pm$ 0.06	22.99 $\pm$ 0.94 5.92 $\pm$ 0.03	21.62 $\pm$ 0.58 4.07 $\pm$ 0.11	13.77 $\pm$ 0.58 1.94 $\pm$ 0.12	131.49 min 16.33 PFLOPs
Grad-Align [1]	54.87 $\pm$ 1.17 22.97 $\pm$ 13.38	31.86 $\pm$ 0.19 11.36 $\pm$ 6.63	31.60 $\pm$ 0.19 9.63 $\pm$ 5.58	31.54 $\pm$ 0.23 9.50 $\pm$ 5.47	26.18 $\pm$ 0.07 6.43 $\pm$ 3.93	229.34 min 16.33 PFLOPs
FGSM-AT [7]	21.36 $\pm$ 8.36 1.45 $\pm$ 0.45	3.04 $\pm$ 0.67 0.96 $\pm$ 0.04	2.38 $\pm$ 0.44 0.85 $\pm$ 0.15	1.88 $\pm$ 0.32 0.82 $\pm$ 0.17	0.19 $\pm$ 0.03 0.51 $\pm$ 0.49	74.17 min 10.89 PFLOPs
Free-AT [25]	58.29 $\pm$ 2.10 20.93 $\pm$ 4.08	30.47 $\pm$ 0.35 11.10 $\pm$ 2.49	30.01 $\pm$ 0.30 9.87 $\pm$ 2.21	29.96 $\pm$ 0.20 9.79 $\pm$ 2.21	24.34 $\pm$ 0.42 5.84 $\pm$ 1.56	71.97 min 10.89 PFLOPs
COAT [14]	<b>67.56 <math>\pm</math> 1.13</b> 65.67 $\pm$ 0.27	24.55 $\pm$ 0.16 10.85 $\pm$ 1.01	23.23 $\pm$ 0.34 6.05 $\pm$ 0.73	22.70 $\pm$ 0.33 4.45 $\pm$ 0.61	18.93 $\pm$ 0.07 2.75 $\pm$ 0.38	88.85 min 13.61 PFLOPs
GAT [26]	65.24 $\pm$ 0.26 <b>72.42 <math>\pm</math> 0.58</b>	27.61 $\pm$ 0.14 3.19 $\pm$ 0.04	26.69 $\pm$ 0.19 1.42 $\pm$ 0.02	26.54 $\pm$ 0.16 0.57 $\pm$ 0.06	21.97 $\pm$ 0.15 0.06 $\pm$ 0.03	126.85 min 16.33 PFLOPs
LIET (Ours)	51.52 $\pm$ 0.38 35.00 $\pm$ 2.85	<b>32.92 <math>\pm</math> 0.12</b> <b>20.34 <math>\pm</math> 0.86</b>	<b>32.75 <math>\pm</math> 0.15</b> <b>17.60 <math>\pm</math> 0.45</b>	<b>32.74 <math>\pm</math> 0.14</b> <b>17.22 <math>\pm</math> 0.42</b>	<b>27.05 <math>\pm</math> 0.09</b> <b>12.04 <math>\pm</math> 0.30</b>	103.74 min 11.11 PFLOPs

Table 7. Comparison of clean accuracy, robust accuracy and training cost (time in minutes and computation in PFLOPs) on CIFAR-100 using ResNet-18. Each method is evaluated with perturbation sizes of 8/255 (first row) and 16/255 (second row). Best results are highlighted in **bold**.

Method	Clean (%) $\uparrow$	PGD-10 (%) $\uparrow$	PGD-20 (%) $\uparrow$	PGD-50 (%) $\uparrow$	AA (%) $\uparrow$	Training Cost
PGD-AT [16]	43.60 $\pm$ 2.45	20.20 $\pm$ 1.82	19.90 $\pm$ 1.41	19.86 $\pm$ 1.23	16.00 $\pm$ 1.02	2432.63 min
	41.27 $\pm$ 2.45	14.85 $\pm$ 0.27	13.44 $\pm$ 0.04	13.15 $\pm$ 0.00	9.54 $\pm$ 0.19	479.01 PFLOPs
PGD-AT-WA [12]	46.23 $\pm$ 0.85	26.09 $\pm$ 0.36	26.06 $\pm$ 0.34	26.06 $\pm$ 0.34	19.62 $\pm$ 0.24	2439.34 min
	41.89 $\pm$ 0.20	19.91 $\pm$ 0.17	18.60 $\pm$ 0.03	18.43 $\pm$ 0.14	12.31 $\pm$ 0.02	479.01 PFLOPs
N-FGSM [3]	47.73 $\pm$ 0.45	25.30 $\pm$ 0.26	25.18 $\pm$ 0.24	25.10 $\pm$ 0.23	18.76 $\pm$ 0.24	498.95 min
	37.98 $\pm$ 1.31	18.11 $\pm$ 0.05	16.94 $\pm$ 0.08	16.64 $\pm$ 0.14	11.11 $\pm$ 0.18	87.09 PFLOPs
FGSM-RS [28]	43.10 $\pm$ 4.09	22.91 $\pm$ 1.35	22.71 $\pm$ 1.30	22.70 $\pm$ 1.27	16.28 $\pm$ 1.03	493.77 min
	5.64 $\pm$ 0.03	2.33 $\pm$ 0.01	2.26 $\pm$ 0.00	2.25 $\pm$ 0.01	1.33 $\pm$ 0.01	87.09 PFLOPs
FGSM-PGI [13]	48.59 $\pm$ 0.19	<b>26.68 <math>\pm</math> 0.25</b>	26.46 $\pm$ 0.15	26.39 $\pm$ 0.17	19.52 $\pm$ 0.07	652.09 min
	24.81 $\pm$ 1.29	5.24 $\pm$ 0.06	4.29 $\pm$ 0.04	4.02 $\pm$ 0.14	1.33 $\pm$ 0.27	87.09 PFLOPs
FGSM-UAP [20]	45.69 $\pm$ 0.99	26.12 $\pm$ 0.05	25.91 $\pm$ 0.02	25.84 $\pm$ 0.04	19.45 $\pm$ 0.14	717.3 min
	16.89 $\pm$ 9.00	2.55 $\pm$ 0.10	2.05 $\pm$ 0.27	1.89 $\pm$ 0.40	0.71 $\pm$ 0.59	130.64 PFLOPs
NuAT [27]	45.55 $\pm$ 0.99	26.51 $\pm$ 0.17	26.38 $\pm$ 0.17	26.37 $\pm$ 0.14	19.55 $\pm$ 0.12	865.64 min
	47.98 $\pm$ 4.95	13.83 $\pm$ 2.69	11.25 $\pm$ 4.00	10.59 $\pm$ 4.58	5.47 $\pm$ 3.37	130.64 PFLOPs
Grad-Align [1]	46.16 $\pm$ 2.03	24.61 $\pm$ 0.51	24.30 $\pm$ 0.40	24.21 $\pm$ 0.41	17.65 $\pm$ 0.28	1514.23 min
	36.83 $\pm$ 2.61	15.57 $\pm$ 0.76	14.43 $\pm$ 0.79	14.22 $\pm$ 0.72	8.71 $\pm$ 0.59	130.64 PFLOPs
FGSM-AT [7]	34.68 $\pm$ 15.85	17.12 $\pm$ 8.69	16.92 $\pm$ 8.60	16.84 $\pm$ 8.55	16.60 $\pm$ 0.36	492.96 min
	5.78 $\pm$ 1.95	1.75 $\pm$ 0.08	1.62 $\pm$ 0.16	1.58 $\pm$ 0.21	0.84 $\pm$ 0.29	87.09 PFLOPs
Free-AT [25]	48.19 $\pm$ 1.78	23.85 $\pm$ 0.24	23.63 $\pm$ 0.25	23.58 $\pm$ 0.26	16.34 $\pm$ 0.32	474.87 min
	32.57 $\pm$ 3.72	14.89 $\pm$ 1.00	14.07 $\pm$ 0.75	13.99 $\pm$ 0.73	8.42 $\pm$ 0.58	87.09 PFLOPs
COAT [14]	<b>59.30 <math>\pm</math> 0.38</b>	18.45 $\pm$ 0.08	17.57 $\pm$ 0.16	17.25 $\pm$ 0.22	11.56 $\pm$ 0.02	578.07 min
	<b>58.52 <math>\pm</math> 1.41</b>	10.65 $\pm$ 0.96	8.15 $\pm$ 0.84	7.38 $\pm$ 0.86	3.71 $\pm$ 0.60	108.87 PFLOPs
GAT [26]	57.68 $\pm$ 0.25	17.97 $\pm$ 0.23	17.26 $\pm$ 0.24	17.01 $\pm$ 0.28	11.53 $\pm$ 0.15	851.30 min
	33.84 $\pm$ 0.91	1.02 $\pm$ 0.07	0.69 $\pm$ 0.04	0.60 $\pm$ 0.02	0.20 $\pm$ 0.01	130.64 PFLOPs
LIET (Ours)	44.73 $\pm$ 0.33	<b>26.68 <math>\pm</math> 0.11</b>	<b>26.54 <math>\pm</math> 0.12</b>	<b>26.54 <math>\pm</math> 0.11</b>	<b>19.79 <math>\pm</math> 0.04</b>	682.47 min
	35.18 $\pm$ 0.76	<b>18.43 <math>\pm</math> 0.13</b>	<b>17.47 <math>\pm</math> 0.10</b>	<b>17.31 <math>\pm</math> 0.19</b>	<b>11.29 <math>\pm</math> 0.15</b>	90.56 PFLOPs

Table 8. Comparison of clean accuracy, robust accuracy and training cost (time in minutes and computation in PFLOPs) on Tiny-ImageNet using PreResNet-18. Each method is evaluated with perturbation sizes of 8/255 (first row) and 12/255 (second row). Best results are highlighted in **bold**.



Method	Clean (%) $\uparrow$	PGD-10 (%) $\uparrow$	PGD-20 (%) $\uparrow$	PGD-50 (%) $\uparrow$	AA (%) $\uparrow$
N-FGSM [3]	80.60 68.10	56.06 39.17	55.76 31.62	55.67 29.86	50.31 21.66
FGSM-RS [28]	89.23 9.99	15.94 9.82	8.94 9.67	5.50 9.56	0.00 7.51
FGSM-PGI [13]	84.25 89.75	<b>60.03</b> 22.92	59.36 18.08	59.23 14.02	52.99 0.30
NuAT [27]	82.68 93.17	57.55 19.34	56.99 10.61	56.94 5.71	52.50 0.12
Grad-Align [1]	82.55 10.00	56.55 9.97	55.99 9.73	55.91 9.12	50.57 2.95
FGSM-AT [7]	<b>89.99</b> 11.46	27.07 9.73	18.51 9.51	11.01 9.34	0.01 7.20
Free-AT [25]	76.98 35.74	51.43 21.77	51.04 20.12	50.93 20.11	46.30 15.48
COAT [14]	86.69 <b>94.51</b>	53.47 34.91	53.08 20.97	53.11 8.29	41.90 0.04
GAT [26]	79.08 93.23	45.45 11.16	45.02 4.59	44.93 1.53	40.39 0.03
LIET (Ours)	82.49 68.32	59.99 <b>42.76</b>	<b>59.53</b> <b>35.96</b>	<b>59.48</b> <b>33.98</b>	<b>53.15</b> <b>23.84</b>

Table 9. Comparison of clean accuracy and robust accuracy on CIFAR-10 using WideResNet-28-10. Each method is evaluated with perturbation sizes of 8/255 (first row) and 16/255 (second row). Best results are highlighted in **bold**.

Method	Clean (%) $\uparrow$	PGD-10 (%) $\uparrow$	PGD-20 (%) $\uparrow$	PGD-50 (%) $\uparrow$	AA (%) $\uparrow$
N-FGSM [3]	57.55 43.49	<b>33.33</b> <b>20.50</b>	32.91 16.54	32.85 15.76	27.17 12.17
FGSM-RS [28]	64.43 2.82	8.31 1.07	5.94 0.91	4.06 0.88	0.00 0.40
FGSM-PGI [13]	<b>70.84</b> 57.42	19.50 9.99	15.89 6.31	13.39 5.29	1.54 3.50
NuAT [27]	66.23 67.24	22.87 12.45	19.50 6.97	17.98 4.41	12.84 1.68
Grad-Align [1]	26.80 1.80	12.25 1.34	12.19 1.13	12.20 1.01	10.00 0.64
FGSM-AT [7]	65.09 1.00	4.23 0.97	1.95 0.60	1.11 0.48	0.00 0.00
Free-AT [25]	49.72 20.94	27.57 10.64	27.49 9.21	27.48 9.28	22.20 5.67
COAT [14]	70.81 <b>77.49</b>	25.97 13.93	23.96 7.98	23.02 3.21	20.16 0.00
GAT [26]	70.66 74.93	27.62 4.30	26.67 1.84	26.62 0.93	23.09 0.04
LIET (Ours)	52.72 35.36	33.30 20.17	<b>33.16</b> <b>17.47</b>	<b>33.04</b> <b>17.19</b>	<b>27.33</b> <b>12.72</b>

Table 10. Comparison of clean accuracy and robust accuracy on CIFAR-100 using WideResNet-28-10. Each method is evaluated with perturbation sizes of 8/255 (first row) and 16/255 (second row). Best results are highlighted in **bold**.