

1 **Organization.** The appendix is organized as follows: In Section A, we describe the hyperparameters  
2 and provide the description about evaluation for multiple perturbations. Furthermore, we provide a  
3 breakdown of all the attacks across various datasets. In addition, we include the previous reviews for  
4 our paper and the changes made to the manuscript.

## 5 A Experimental setup

### 6 A.1 Datasets

- 7 1. **CIFAR-10.** This dataset [1] contains 60,000 images with 5,000 images for training and 1,000  
8 images for test for each class. Each image is sized  $32 \times 32$ , we use the Wide ResNet 28-10  
9 architecture [2] as a base network for this dataset.
- 10 2. **SVHN.** This dataset [3] contains 73257 training and 26032 testing images of digits and numbers  
11 in natural scene images containing ten-digit classes. Each image is sized  $32 \times 32$ , we use the Wide  
12 ResNet 28-10 architecture similar to the CIFAR-10 dataset as the base network.
- 13 3. **Tiny-ImageNet.** This dataset <sup>1</sup> is a subset of ImageNet [4] dataset, consisting of 500, 50, and 50  
14 images for training, validation, and test dataset, respectively. This dataset contains  $64 \times 64$  size  
15 images from 200 classes, we use ResNet50 [5] as a base network for this dataset.

### 16 A.2 Training setup

17 We use the SGD optimizer with momentum 0.9 and weight decay  $5 \cdot 10^{-4}$  to train all our models  
18 with cyclic learning rate with a maximum learning rate  $\lambda$  that increases linearly from 0 to  $\lambda$  over first  
19  $N/2$  epochs and then decreases linearly from  $N/2$  to 0 in the remainder epochs, as recommended by  
20 [6] for fast convergence of adversarial training. We train all the models for 30 epochs on a single  
21 machine with four GeForce RTX 2080Ti using WideResNet 28-10 architecture [2]. We use the  
22 maximum learning rate of  $\lambda = 0.21$  for all our experiments. We use  $\beta = 16$  for all the experiments  
23 with our meta noise generator. The generator is formulated as a convolutional network with four  
24  $3 \times 3$  convolutional layers with LeakyReLU activations and one residual connection from input to  
25 output. All our algorithms are implemented in Pytorch [7]. We use the weight for the KL divergence  
26 ( $\beta = 6.0$ ) for TRADES and RST in all our experiments. We replicate all the baselines on SVHN  
27 and TinyImageNet since most of the baseline methods have reported their results on MNIST and  
28 CIFAR-10. Unfortunately, we found that MSD [8] did not converge for larger datasets even after our  
29 extensive hyperparameter-search. We believe that this is due to the change in formulation of the  
30 inner optimization which leads to a difficulty in convergence for larger datasets. Since the authors  
31 also report their results on CIFAR-10, we do not use it as a baseline for other datasets.

### 32 A.3 Evaluation setup

33 For  $\ell_\infty$  perturbations, we use PGD [9], Brendel and Bethge attack [10], and AutoAttack [11]. For  
34  $\ell_2$  perturbations, we use CarliniWagner attack [12], PGD [9], Brendel and Bethge attack [10], and  
35 AutoAttack [11]. For  $\ell_1$  perturbations, we use SLIDE [13], Salt and pepper [14], and EAD attack [15].  
36 For all our experiments and evaluation, we use  $\varepsilon = \{0.03, 8, 0.31\}$  and  $\alpha = \{0.004, 1.0, 0.1\}$   
37 for  $\ell_\infty, \ell_1$ , and  $\ell_2$  attacks for CIFAR-10 and SVHN respectively. For Tiny-ImageNet we use  
38  $\varepsilon = \{0.01, 8, 0.31\}$  and  $\alpha = \{0.004, 1.0, 0.1\}$  for  $\ell_\infty, \ell_1$ , and  $\ell_2$  attacks respectively. We use 10  
39 steps of PGD attack for  $\ell_\infty, \ell_2$  during training. For  $\ell_1$  adversarial training, we use 20 steps during  
40 training and 100 steps during evaluation. We use the code provided by the authors for evaluation  
41 against AutoAttack [11] and Foolbox [14] library for all the other attacks.

42 Due to the length limit of our paper, we provide a breakdown of all the attacks on CIFAR-10 in  
43 Table A.1, SVHN on Wide ResNet 28-10 in Table A.2, Tiny-ImageNet on ResNet50 in Table A.3.  
44 Besides, we analyze the noise learned by our meta-learning framework on multiple datasets and the  
45 loss landscape on the CIFAR-10 dataset.

---

<sup>1</sup><https://tiny-imagenet.herokuapp.com/>

Table A.1: Summary of adversarial accuracy results for CIFAR-10 on Wide ResNet 28-10 architecture.

	$\text{Adv}_\infty$	$\text{Adv}_1$	$\text{Adv}_2$	$\text{Trades}_\infty$	$\text{RST}_\infty$	$\text{Adv}_{\text{avg}}$	$\text{Adv}_{\text{max}}$	MSD	MNG-AC
Clean Accuracy	86.8 $\pm$ 0.1	93.3 $\pm$ 0.6	91.7 $\pm$ 0.2	84.7 $\pm$ 0.3	88.9 $\pm$ 0.2	87.1 $\pm$ 0.2	85.4 $\pm$ 0.3	82.3 $\pm$ 0.2	84.9 $\pm$ 0.3
PGD- $\ell_\infty$	46.9 $\pm$ 0.5	0.40 $\pm$ 0.7	23.6 $\pm$ 0.2	52.0 $\pm$ 0.6	56.9 $\pm$ 0.1	35.2 $\pm$ 0.8	42.2 $\pm$ 1.1	45.4 $\pm$ 0.4	44.5 $\pm$ 1.1
PGD-Foolbox	54.7 $\pm$ 0.4	0.33 $\pm$ 0.6	35.3 $\pm$ 0.4	57.8 $\pm$ 0.5	62.9 $\pm$ 0.3	45.0 $\pm$ 0.4	50.4 $\pm$ 0.4	51.7 $\pm$ 0.8	50.8 $\pm$ 0.8
AutoAttack	44.9 $\pm$ 0.7	0.0 $\pm$ 0.0	20.7 $\pm$ 0.4	48.8 $\pm$ 1.1	53.9 $\pm$ 0.3	33.8 $\pm$ 0.7	39.9 $\pm$ 0.9	42.7 $\pm$ 0.2	42.8 $\pm$ 0.8
Brendel & Bethge	49.9 $\pm$ 1.1	0.0 $\pm$ 0.0	26.8 $\pm$ 0.3	52.1 $\pm$ 0.7	56.5 $\pm$ 1.8	39.6 $\pm$ 0.7	45.8 $\pm$ 0.9	48.3 $\pm$ 0.4	46.8 $\pm$ 0.9
<b>All <math>\ell_\infty</math> attacks</b>	44.9 $\pm$ 0.7	0.0 $\pm$ 0.0	20.7 $\pm$ 0.3	48.9 $\pm$ 0.7	54.9 $\pm$ 1.8	33.8 $\pm$ 0.7	39.9 $\pm$ 0.9	43.7 $\pm$ 0.2	42.2 $\pm$ 0.9
PGD- $\ell_1$	12.8 $\pm$ 0.6	91.6 $\pm$ 1.4	27.7 $\pm$ 0.7	17.9 $\pm$ 0.6	22.0 $\pm$ 0.5	49.0 $\pm$ 0.3	44.6 $\pm$ 0.2	46.8 $\pm$ 1.4	55.0 $\pm$ 1.2
PGD-Foolbox	35.2 $\pm$ 0.7	92.3 $\pm$ 1.3	53.1 $\pm$ 0.5	40.3 $\pm$ 0.7	44.6 $\pm$ 0.3	64.5 $\pm$ 0.2	60.7 $\pm$ 0.5	60.3 $\pm$ 0.4	65.5 $\pm$ 0.1
EAD	72.9 $\pm$ 1.0	87.1 $\pm$ 3.3	75.9 $\pm$ 1.9	80.2 $\pm$ 0.7	84.5 $\pm$ 0.2	85.7 $\pm$ 0.2	83.3 $\pm$ 0.5	80.8 $\pm$ 0.1	79.3 $\pm$ 0.6
SAPA	71.5 $\pm$ 0.2	80.2 $\pm$ 1.8	81.9 $\pm$ 0.5	71.4 $\pm$ 0.7	76.0 $\pm$ 0.5	82.7 $\pm$ 0.1	80.0 $\pm$ 0.1	76.9 $\pm$ 0.5	76.7 $\pm$ 0.4
<b>All <math>\ell_1</math> attacks</b>	12.8 $\pm$ 0.6	78.1 $\pm$ 1.8	27.7 $\pm$ 0.7	17.9 $\pm$ 0.6	22.0 $\pm$ 0.5	49.0 $\pm$ 0.3	44.6 $\pm$ 0.2	43.7 $\pm$ 0.2	55.0 $\pm$ 1.2
PGD- $\ell_2$	78.7 $\pm$ 0.3	47.6 $\pm$ 1.6	84.6 $\pm$ 0.2	77.0 $\pm$ 0.9	82.2 $\pm$ 0.2	81.5 $\pm$ 0.2	79.1 $\pm$ 0.3	76.5 $\pm$ 0.1	75.6 $\pm$ 0.4
PGD-Foolbox	74.6 $\pm$ 0.2	5.1 $\pm$ 2.1	79.8 $\pm$ 0.2	73.3 $\pm$ 0.6	78.3 $\pm$ 0.2	77.6 $\pm$ 0.2	75.8 $\pm$ 0.3	73.6 $\pm$ 0.5	73.4 $\pm$ 0.1
Gaussian Noise	85.2 $\pm$ 0.4	88.5 $\pm$ 1.8	90.5 $\pm$ 1.1	83.2 $\pm$ 0.3	87.8 $\pm$ 0.2	86.2 $\pm$ 0.5	83.3 $\pm$ 0.3	70.9 $\pm$ 1.1	79.3 $\pm$ 0.1
AutoAttack	69.9 $\pm$ 0.4	0.0 $\pm$ 0.0	76.8 $\pm$ 0.4	69.4 $\pm$ 0.3	73.7 $\pm$ 0.1	74.9 $\pm$ 0.4	73.2 $\pm$ 0.2	71.9 $\pm$ 0.4	71.5 $\pm$ 0.1
Brendel & Bethge	71.8 $\pm$ 0.9	0.0 $\pm$ 0.0	78.1 $\pm$ 0.6	70.2 $\pm$ 0.1	75.0 $\pm$ 0.3	75.9 $\pm$ 0.3	74.1 $\pm$ 0.4	80.4 $\pm$ 0.4	72.3 $\pm$ 0.1
CWL2	70.5 $\pm$ 0.2	0.1 $\pm$ 0.0	77.2 $\pm$ 0.5	69.7 $\pm$ 0.3	74.2 $\pm$ 0.1	74.6 $\pm$ 1.2	73.5 $\pm$ 0.2	71.1 $\pm$ 1.1	71.0 $\pm$ 0.1
<b>All <math>\ell_2</math> attacks</b>	69.3 $\pm$ 0.4	0.0 $\pm$ 0.0	76.8 $\pm$ 0.4	69.4 $\pm$ 0.3	73.6 $\pm$ 0.1	74.9 $\pm$ 0.4	73.2 $\pm$ 0.2	70.6 $\pm$ 1.1	71.5 $\pm$ 0.1
$\text{Acc}_{\text{adv}}^{\text{union}}$	12.9 $\pm$ 0.5	0.0 $\pm$ 0.0	17.9 $\pm$ 0.8	17.2 $\pm$ 0.6	21.1 $\pm$ 1.0	31.0 $\pm$ 1.4	35.7 $\pm$ 0.3	35.8 $\pm$ 0.1	<b>41.6<math>\pm</math>0.8</b>
$\text{Acc}_{\text{adv}}^{\text{avg}}$	42.6 $\pm$ 0.4	25.1 $\pm$ 1.6	47.6 $\pm$ 0.4	45.4 $\pm$ 0.3	50.2 $\pm$ 0.5	52.6 $\pm$ 0.5	52.5 $\pm$ 0.3	52.0 $\pm$ 0.4	<b>56.2<math>\pm</math>0.2</b>

Table A.2: Summary of adversarial accuracy results for SVHN dataset on Wide ResNet 28-10 architecture.

	$\text{Adv}_\infty$	$\text{Adv}_1$	$\text{Adv}_2$	$\text{Trades}_\infty$	$\text{RST}_\infty$	$\text{Adv}_{\text{avg}}$	$\text{Adv}_{\text{max}}$	MNG-AC
Clean Accuracy	92.8 $\pm$ 0.1	92.4 $\pm$ 1.6	94.9 $\pm$ 0.0	93.9 $\pm$ 0.0	95.6 $\pm$ 0.0	92.6 $\pm$ 0.1	88.2 $\pm$ 1.6	93.4 $\pm$ 0.0
PGD- $\ell_\infty$	49.1 $\pm$ 0.1	3.2 $\pm$ 2.4	29.4 $\pm$ 0.1	55.5 $\pm$ 1.4	66.9 $\pm$ 0.8	22.4 $\pm$ 3.1	36.6 $\pm$ 2.0	40.5 $\pm$ 0.1
PGD-Foolbox	60.7 $\pm$ 0.4	2.5 $\pm$ 1.9	47.6 $\pm$ 0.6	66.4 $\pm$ 1.1	73.8 $\pm$ 0.3	32.5 $\pm$ 3.2	49.9 $\pm$ 0.0	57.5 $\pm$ 1.8
AutoAttack	46.2 $\pm$ 0.6	0.0 $\pm$ 0.0	18.9 $\pm$ 0.5	49.9 $\pm$ 1.8	61.0 $\pm$ 2.0	17.6 $\pm$ 2.6	17.5 $\pm$ 0.9	33.7 $\pm$ 0.0
Brendel & Bethge	51.6 $\pm$ 0.7	0.0 $\pm$ 0.0	22.9 $\pm$ 0.8	55.8 $\pm$ 1.5	65.6 $\pm$ 1.2	20.2 $\pm$ 2.9	6.3 $\pm$ 2.3	40.0 $\pm$ 0.3
<b>All <math>\ell_\infty</math> attacks</b>	46.2 $\pm$ 0.6	0.0 $\pm$ 0.0	18.7 $\pm$ 0.6	49.9 $\pm$ 1.7	60.9 $\pm$ 2.0	17.4 $\pm$ 2.3	5.9 $\pm$ 1.2	35.1 $\pm$ 1.9
PGD- $\ell_1$	3.1 $\pm$ 0.3	95.0 $\pm$ 1.8	30.5 $\pm$ 0.4	1.7 $\pm$ 0.3	0.7 $\pm$ 0.6	55.8 $\pm$ 2.1	48.4 $\pm$ 2.9	44.5 $\pm$ 3.2
PGD-Foolbox	19.9 $\pm$ 0.8	94.6 $\pm$ 0.4	57.5 $\pm$ 0.1	15.5 $\pm$ 0.2	11.3 $\pm$ 0.5	79.2 $\pm$ 3.4	85.4 $\pm$ 3.2	75.2 $\pm$ 2.8
EAD	65.7 $\pm$ 2.1	87.8 $\pm$ 1.9	82.3 $\pm$ 1.2	51.5 $\pm$ 2.9	60.4 $\pm$ 0.8	84.8 $\pm$ 2.4	84.5 $\pm$ 3.8	86.2 $\pm$ 2.2
SAPA	79.4 $\pm$ 0.8	77.3 $\pm$ 5.2	87.3 $\pm$ 0.1	73.5 $\pm$ 1.0	86.2 $\pm$ 0.5	88.5 $\pm$ 0.6	80.9 $\pm$ 4.0	89.9 $\pm$ 1.6
<b>All <math>\ell_1</math> attacks</b>	3.0 $\pm$ 0.3	77.9 $\pm$ 6.3	30.3 $\pm$ 0.3	1.6 $\pm$ 0.3	0.7 $\pm$ 0.6	54.2 $\pm$ 2.9	48.3 $\pm$ 4.1	47.4 $\pm$ 2.2
PGD- $\ell_2$	81.6 $\pm$ 0.5	3.9 $\pm$ 1.4	87.8 $\pm$ 0.2	83.9 $\pm$ 0.8	85.3 $\pm$ 0.2	85.6 $\pm$ 0.6	84.3 $\pm$ 1.1	90.4 $\pm$ 0.6
PGD-Foolbox	73.2 $\pm$ 0.2	1.9 $\pm$ 1.8	82.8 $\pm$ 0.6	75.0 $\pm$ 0.7	76.0 $\pm$ 0.3	80.6 $\pm$ 0.1	60.1 $\pm$ 0.8	86.1 $\pm$ 0.1
Gaussian Noise	92.1 $\pm$ 0.2	16.5 $\pm$ 4.2	94.2 $\pm$ 0.2	93.3 $\pm$ 1.4	94.2 $\pm$ 0.6	92.2 $\pm$ 0.2	83.8 $\pm$ 0.6	93.2 $\pm$ 0.4
AutoAttack	59.0 $\pm$ 0.7	0.0 $\pm$ 0.0	79.3 $\pm$ 0.1	56.4 $\pm$ 1.3	60.7 $\pm$ 0.6	75.6 $\pm$ 0.1	40.0 $\pm$ 2.3	78.0 $\pm$ 0.8
Brendel & Bethge	68.2 $\pm$ 0.5	0.0 $\pm$ 0.0	81.0 $\pm$ 0.1	64.8 $\pm$ 0.9	68.1 $\pm$ 0.5	76.4 $\pm$ 0.4	32.7 $\pm$ 3.8	78.4 $\pm$ 0.4
CWL2	63.5 $\pm$ 0.8	0.1 $\pm$ 0.1	80.1 $\pm$ 1.4	61.4 $\pm$ 0.3	63.9 $\pm$ 0.2	76.8 $\pm$ 0.1	55.3 $\pm$ 5.2	80.9 $\pm$ 0.9
<b>All <math>\ell_2</math> attacks</b>	59.2 $\pm$ 0.7	0.0 $\pm$ 0.0	79.3 $\pm$ 0.1	56.0 $\pm$ 1.4	60.6 $\pm$ 0.6	74.7 $\pm$ 0.1	31.0 $\pm$ 5.0	77.6 $\pm$ 1.0
$\text{Acc}_{\text{adv}}^{\text{union}}$	3.0 $\pm$ 0.3	0.0 $\pm$ 0.0	16.4 $\pm$ 0.7	1.6 $\pm$ 0.3	0.7 $\pm$ 0.6	16.6 $\pm$ 1.3	5.8 $\pm$ 1.7	<b>30.3<math>\pm</math>1.8</b>
$\text{Acc}_{\text{adv}}^{\text{avg}}$	36.2 $\pm$ 0.3	23.9 $\pm$ 2.1	42.8 $\pm$ 0.2	35.8 $\pm$ 0.6	40.7 $\pm$ 0.8	43.0 $\pm$ 1.0	26.7 $\pm$ 2.5	<b>52.6<math>\pm</math>0.5</b>

Table A.3: Summary of adversarial accuracy results for Tiny-ImageNet on ResNet50 architecture.

	Adv $_{\infty}$	Adv $_1$	Adv $_2$	Trades $_{\infty}$	Adv $_{\text{avg}}$	Adv $_{\text{max}}$	MNG-AC
Clean Accuracy	54.2 $\pm$ 0.1	57.8 $\pm$ 0.2	59.8 $\pm$ 0.1	48.2 $\pm$ 0.2	56.0 $\pm$ 0.2	53.5 $\pm$ 0.0	53.1 $\pm$ 0.1
PGD- $\ell_{\infty}$	32.1 $\pm$ 0.0	11.5 $\pm$ 1.2	17.9 $\pm$ 1.1	32.2 $\pm$ 0.4	25.0 $\pm$ 0.6	32.0 $\pm$ 0.6	29.3 $\pm$ 0.3
PGD-Foolbox	34.6 $\pm$ 0.4	17.2 $\pm$ 0.1	5.2 $\pm$ 0.6	34.1 $\pm$ 0.2	34.0 $\pm$ 0.2	28.3 $\pm$ 0.1	32.3 $\pm$ 0.3
AutoAttack	29.6 $\pm$ 0.1	10.1 $\pm$ 0.7	16.3 $\pm$ 0.3	28.7 $\pm$ 0.9	23.7 $\pm$ 0.2	30.0 $\pm$ 0.1	27.7 $\pm$ 0.4
Brendel & Bethge	32.7 $\pm$ 0.1	14.6 $\pm$ 0.8	20.8 $\pm$ 0.6	31.0 $\pm$ 0.9	28.1 $\pm$ 0.2	33.2 $\pm$ 0.5	31.5 $\pm$ 0.6
<b>All <math>\ell_{\infty}</math> attacks</b>	29.6 $\pm$ 0.1	10.5 $\pm$ 0.7	5.2 $\pm$ 0.6	28.7 $\pm$ 0.9	23.7 $\pm$ 0.2	29.8 $\pm$ 0.1	27.4 $\pm$ 0.7
PGD- $\ell_1$	32.0 $\pm$ 1.1	39.3 $\pm$ 0.9	37.2 $\pm$ 0.2	31.1 $\pm$ 0.3	38.0 $\pm$ 0.1	33.6 $\pm$ 0.4	39.0 $\pm$ 0.9
PGD-Foolbox	40.0 $\pm$ 0.8	44.8 $\pm$ 0.2	45.2 $\pm$ 0.2	37.6 $\pm$ 0.9	44.7 $\pm$ 1.5	40.6 $\pm$ 0.1	45.0 $\pm$ 0.2
EAD	52.3 $\pm$ 1.5	56.3 $\pm$ 0.6	57.3 $\pm$ 0.0	46.7 $\pm$ 0.9	54.6 $\pm$ 0.9	51.2 $\pm$ 0.2	52.7 $\pm$ 0.3
SAPA	46.5 $\pm$ 0.9	52.9 $\pm$ 0.7	53.5 $\pm$ 1.2	40.8 $\pm$ 0.1	50.3 $\pm$ 1.1	46.6 $\pm$ 0.1	49.3 $\pm$ 0.4
<b>All <math>\ell_1</math> attacks</b>	31.8 $\pm$ 1.0	39.3 $\pm$ 1.0	37.2 $\pm$ 0.4	30.9 $\pm$ 0.2	38.0 $\pm$ 0.2	33.4 $\pm$ 0.3	39.6 $\pm$ 0.7
PGD- $\ell_2$	48.5 $\pm$ 1.1	49.1 $\pm$ 0.1	51.8 $\pm$ 1.8	42.6 $\pm$ 0.7	49.9 $\pm$ 1.7	47.0 $\pm$ 0.3	49.1 $\pm$ 0.4
PGD-Foolbox	45.6 $\pm$ 0.4	45.2 $\pm$ 0.4	47.7 $\pm$ 0.7	41.0 $\pm$ 0.3	47.0 $\pm$ 1.3	44.9 $\pm$ 0.4	47.0 $\pm$ 0.2
Gaussian Noise	52.5 $\pm$ 1.3	56.1 $\pm$ 0.6	57.6 $\pm$ 0.3	46.4 $\pm$ 0.9	54.4 $\pm$ 0.8	51.1 $\pm$ 0.0	52.1 $\pm$ 0.5
AutoAttack	42.4 $\pm$ 0.8	41.9 $\pm$ 0.0	44.6 $\pm$ 0.6	38.9 $\pm$ 0.8	44.4 $\pm$ 1.3	42.4 $\pm$ 0.9	44.6 $\pm$ 0.4
Brendel & Bethge	43.7 $\pm$ 0.4	44.4 $\pm$ 0.1	46.6 $\pm$ 1.1	39.2 $\pm$ 0.7	45.1 $\pm$ 1.6	43.6 $\pm$ 0.4	45.4 $\pm$ 0.1
CWL2	43.5 $\pm$ 1.3	44.8 $\pm$ 1.1	47.5 $\pm$ 0.7	39.5 $\pm$ 0.4	46.8 $\pm$ 1.9	43.4 $\pm$ 0.1	46.0 $\pm$ 0.4
<b>All <math>\ell_2</math> attacks</b>	42.5 $\pm$ 0.6	41.9 $\pm$ 0.0	44.9 $\pm$ 0.1	35.8 $\pm$ 0.7	44.6 $\pm$ 0.1	42.4 $\pm$ 1.0	44.8 $\pm$ 0.1
Acc $_{\text{adv}}^{\text{union}}$	19.8 $\pm$ 1.1	10.1 $\pm$ 0.7	5.2 $\pm$ 0.6	26.1 $\pm$ 0.9	23.6 $\pm$ 0.3	<b>29.0<math>\pm</math> 0.3</b>	27.4 $\pm$ 0.8
Acc $_{\text{adv}}^{\text{avg}}$	33.8 $\pm$ 0.1	30.4 $\pm$ 0.1	29.1 $\pm$ 0.0	32.8 $\pm$ 0.1	35.4 $\pm$ 0.7	35.3 $\pm$ 0.4	<b>37.2<math>\pm</math> 0.6</b>

---

### Summary of previous reviews

---

All the reviewers found our work novel [R2, R3] and tackles an important problem of robustness against multiple perturbations [R2, R3, R4]. Moreover, all reviewers appreciated the significant decrease in computation time for multiple perturbations adversarial training. However, we guess that the notations were unclear, and some crucial experiments were missing in the paper. The paper has now been rewritten, provides exhaustive experiments on benchmark datasets against state-of-the-art adversarial attacks.

**Ablation studies** [R1, R4]. Following R1, R4 suggestions, we now provide a comparison with Gaussian Noise, MNG without the AC loss on multiple datasets in our ablation studies. We can observe that MNG-AC significantly outperform all the baselines, which is not surprising since learning the noise is more efficient way to improve generalization.

**Comparison with baselines** [R1, R2]. R1, R2 suggests to compare against Maini et al. [8], however, since it was accepted to ICML 20, we did not have sufficient time to add it as a baseline for the NeurIPS submission. We have now updated our manuscript and provide a comparison with it in the main paper.

**Stronger attacks** [R1, R3]. Following R1 and R3's suggestion, we show a comparison against Autoattack [11] on all the datasets. Furthermore, we compare against Brendel and Bethge attack [10] for various  $\ell_p$  norms.

**Scalability to larger datasets** [R3, R4]. Both the reviewers mention that our paper could be made stronger with experiments on the larger datasets. We have now updated our paper and provided experiments on large scale architectures and datasets. Due to the limited time provided for rebuttal, we could not perform ImageNet-scale experiments, but we will perform this experiment and include it in the final version if the paper is accepted.

**Incorrect implementation** [R2]. R2 raises an issue with the implementation of our max strategy, we have fixed this in our updated manuscript.

**Other comments.** Reviewers also highlight some inconsistencies with our notations. We have revised our manuscript by clarifying all the notations and inconsistencies.

---

### The changes that have been made to the paper

---

- We rewrote and revised throughout the paper including texts, notations, and figures to avoid confusion and for readability and clarity.
- We included the experimental results using modern deep neural networks, ResNet18 on CIFAR-10 and SVHN -> Wide ResNet 28-10 on CIFAR-10, SVHN and ResNet50 on Tiny-ImageNet. (From [R3, R4] comments)
- We provide a standard deviation for all the results. (From [R1, R4] comments)
- We include additional state-of-the-art baselines. (From [R1, R2] comments)
- We include ablation studies dissecting various components of our framework. (From [R1, R4] comments)
- We evaluate on state-of-the-art adversarial attacks. (From [R1, R3] comments)

## References

- [1] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [2] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference*, 2016.
- [3] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *Workshop on Deep Learning and Unsupervised Feature Learning, NeurIPS*, 2011.
- [4] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 2015.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, 2016.
- [6] Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training. In *ICLR*, 2020.
- [7] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019.
- [8] Pratyush Maini, Eric Wong, and J Zico Kolter. Adversarial robustness against the union of multiple perturbation models. In *ICML*, 2020.
- [9] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2017.
- [10] Wieland Brendel, Jonas Rauber, Matthias Kümmeler, Ivan Ustyuzhaninov, and Matthias Bethge. Accurate, reliable and fast robustness evaluation. In *NeurIPS*, 2019.
- [11] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *ICML*, 2020.
- [12] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, 2017.
- [13] Florian Tramèr and Dan Boneh. Adversarial training and robustness for multiple perturbations. In *NeurIPS*, 2019.
- [14] Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models. In *Reliable Machine Learning in the Wild Workshop, ICML*, 2017.
- [15] Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: Elastic-net attacks to deep neural networks via adversarial examples. In *AAAI*, 2018.