

# Learnable Boundary Guided Adversarial Training

## Supplementary Material

### A. Robustness under Black-box attack

Table 7: Comparison of our method with previous defense models under black-box attack on CIFAR-100 and CIFAR-10. To rule out randomness, the numbers are averaged over 2 independently trained models.  $Acc_n$  represents accuracy on natural images.  $BAcc_r$  represents robustness under black-box attack.  $WAcc_r$  represents robustness under white-box attack

Target Models	$BAcc_r$	$WAcc_r$	$Acc_n$	Source Models	Dataset
TRADES ( $\alpha = 1$ )	61.29%	25.31%	62.37%	Natural	CIFAR-100
TRADES ( $\alpha = 6$ )	55.52%	30.93%	56.51%	Natural	CIFAR-100
LBGAT+ALP	61.38%	<b>35.25%</b>	62.67%	Natural	CIFAR100
LBGAT+TRADES ( $\alpha=0$ )	<b>68.35%</b>	<b>33.01%</b>	<b>70.03%</b>	Natural	CIFAR-100
TRADES ( $\alpha = 1$ )	42.32%	25.31%	62.37%	LBGAT+TRADES ( $\alpha=0$ )	CIFAR-100
TRADES ( $\alpha = 6$ )	41.67%	30.93%	56.51%	LBGAT+TRADES ( $\alpha=0$ )	CIFAR-100
LBGAT+ALP	45.68%	<b>35.25%</b>	62.67%	TRADES ( $\alpha = 6$ )	CIFAR-100
LBGAT+TRADES ( $\alpha=0$ )	<b>50.27%</b>	<b>33.01%</b>	<b>70.03%</b>	TRADES ( $\alpha = 6$ )	CIFAR-100
TRADES ( $\alpha = 1$ )	87.00%	49.14%	<b>88.64%</b>	Natural	CIFAR-10
TRADES ( $\alpha = 6$ )	83.30%	56.61%	84.92%	Natural	CIFAR-10
LBGAT+TRADES ( $\alpha = 0$ )	<b>87.20%</b>	<b>57.55%</b>	<b>88.22%</b>	Natural	CIFAR-10
TRADES ( $\alpha = 1$ )	66.18%	49.14%	<b>88.64%</b>	LBGAT+TRADES( $\alpha=0$ )	CIFAR-10
TRADES ( $\alpha = 6$ )	67.18%	56.61%	84.92%	LBGAT+TRADES ( $\alpha=0$ )	CIFAR-10
LBGAT+TRADES ( $\alpha = 0$ )	<b>68.45%</b>	<b>57.55%</b>	<b>88.22%</b>	TRADES ( $\alpha=6$ )	CIFAR-10

### B. Our Method Creates New SOTA Under the Strongest Auto-Attack on CIFAR-100

To further show the effectiveness of our method, we compare with more previous works. The experimental results are shown in Table 8. On the more challenging CIFAR-100 dataset, our method creates a new state-of-the-art (SOTA) on both robustness and natural accuracy. Specifically, our LBGAT ( $\alpha = 0$ ) model with WideResNet-34-10 architecture significantly outperforms previous SOAT method [6] by 7.08% in the aspect of performance on natural data. Meanwhile, our method surpasses it with respect to model robustness. Further, our strongest model LBGAT ( $\alpha = 6$ ) with WideResNet-34-10 architecture enjoys 2.4% higher robustness than [6].

Moreover, It is worthy to note that our LBGAT ( $\alpha = 6$ ) model achieves even strong robustness than the model, by Hendrycks *et al.* [18], pre-trained on full ImageNet. At the same time, we also surpasses it in the aspect of natural accuracy.

Table 8: More comparisons under the strongest Auto-Attack on CIFAR-100 dataset. "†" denotes numbers are directly copied from [10]. "\*" denotes that the method has used additional unlabeled data.

Methods	Model	Acc <sub>n</sub>	Acc <sub>r</sub>
LBGAT ( $\alpha = 0$ ) Ours	WideResNet-34-20	<b>71.00%</b>	<b>27.66%</b>
LBGAT ( $\alpha = 6$ ) Ours	WideResNet-34-20	62.55%	<b>30.20%</b>
LBGAT ( $\alpha = 0$ ) Ours	WideResNet-34-10	<b>70.03%</b>	<b>27.05%</b>
LBGAT ( $\alpha = 6$ ) Ours	WideResNet-34-10	60.43%	29.34%
TRADES ( $\alpha = 1$ ) [56]	WideResNet-34-10	62.37%	22.24%
TRADES ( $\alpha = 6$ ) [56]	WideResNet-34-10	56.50%	26.87%
Sitawarin <i>et al.</i> [38] †	WideResNet-34-10	62.82%	24.57%
Chen <i>et al.</i> [6] †	WideResNet-34-10	62.15%	26.94%
Hendrycks <i>et al.</i> [18] †*	WideResNet-28-10	59.23%	28.42%
Rice <i>et al.</i> [33] †	ResNet-18	53.83%	18.95%

### C. More Comparisons Under the Strongest Auto-Attack on CIFAR-10

We also compare with more previous methods on CIFAR-10 dataset. The experimental results are summarized in Table 9. Our LBGAT ( $\alpha = 0$ ) model with WideResNet-34-10 architecture can consistently enjoy higher natural performance while keeping the strongest robustness. We observe that though many fast adversarial training methods, like [45, 35] are proposed to accelerate the training process, their performance are usually unsatisfied.

Table 9: More comparisons under the strongest Auto-Attack on CIFAR-10 dataset. "†" denotes numbers are directly copied from [10]. "\*" denotes the methods aiming to accelerate adversarial training.

Methods	Model	Acc <sub>n</sub>	Acc <sub>r</sub>
LBGAT ( $\alpha = 0$ ) Ours	WideResNet-34-20	<b>88.70%</b>	<b>53.58%</b>
LBGAT ( $\alpha = 6$ ) Ours	WideResNet-34-20	83.61%	<b>54.45%</b>
LBGAT ( $\alpha = 0$ ) Ours	WideResNet-34-10	<b>88.22%</b>	<b>52.86%</b>
LBGAT ( $\alpha = 6$ ) Ours	WideResNet-34-10	81.98%	53.14%
Rice <i>et al.</i> [33] †	WideResNet-34-20	85.34%	53.42%
TRADES ( $\alpha = 1$ )	WideResNet-34-10	<b>88.64%</b>	48.11%
TRADES ( $\alpha = 6$ )	WideResNet-34-10	84.92%	52.64%
Kumari <i>et al.</i> [22] †	WideResNet-34-10	87.80%	49.12%
Mao <i>et al.</i> [28] †	WideResNet-34-10	86.21%	47.41%
Zhang <i>et al.</i> [55] †*	WideResNet-34-10	87.20%	44.83%
Shafahi <i>et al.</i> [35] †*	WideResNet-34-10	86.11%	41.47%
Chan <i>et al.</i> [5] †	WideResNet-34-10	<b>93.79%</b>	0.26%
Wang <i>et al.</i> [45] †*	WideResNet-28-10	<b>92.80%</b>	29.35%
Qin <i>et al.</i> [32] †	WideResNet-40-8	86.28%	52.81%
Chen <i>et al.</i> [8] †	ResNet-50	86.04%	51.56%
Xiao <i>et al.</i> [47] †	DenseNet-121	79.28%	18.50%
Wong <i>et al.</i> [46] †	ResNet-18	83.34%	43.21%