

# LAS-AT: Adversarial Training with Learnable Attack Strategy

## Supplementary Material

Xiaojun Jia<sup>1,2,†\*</sup>, Yong Zhang<sup>3,\*</sup>, Baoyuan Wu<sup>4,5,‡</sup>, Ke Ma<sup>6</sup>, Jue Wang<sup>3</sup>, Xiaochun Cao<sup>1,2,‡</sup>

<sup>1</sup>Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

<sup>2</sup>School of Cyberspace Security, University of Chinese Academy of Sciences, Beijing, China

<sup>3</sup>Tencent, AI Lab, Shenzhen, China

<sup>4</sup>School of Data Science, The Chinese University of Hong Kong, Shenzhen, China

<sup>5</sup>Secure Computing Lab of Big Data, Shenzhen Research Institute of Big Data, Shenzhen, China

<sup>6</sup>School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing, China

jiaxiaojun@iie.ac.cn; zhangyong201303@gmail.com; wubaoyuan@cuhk.edu.cn;

make@ucas.ac.cn; arphid@gmail.com; caoxiaochun@iie.ac.cn

In the manuscript, we report extensive experimental results on three benchmark databases. In this supplementary material, more additional studies are provided as follows:

- We present the architecture of the strategy network (see Sec. 1).
- We present the proof of theorem 1 (see Sec. 2).
- We introduce the details of these databases (see Sec. 3).
- We introduce the detailed settings for the proposed LAT-AT (see Sec. 4).
- We introduce the training and evaluation settings for the comparison methods (see Sec. 5).
- We present comparisons with AWP trained with more iterations and a larger perturbation strength (see Sec. 6).
- We study how the maximal perturbation strength, the number of iteration, and the step size affect the robustness and the clean accuracy of the target network in adversarial training. We use PGD-AT [6] as an illustration (see Sec. 7).
- We present the selection of the two trade-off hyper-parameters (*i.e.*,  $\alpha$  and  $\beta$ ) in the objective function (see Sec. 8).

- We illustrate the evolution of the generated perturbation strength of several images during the training process (see Sec. 9).
- We conduct experiments on more image databases (see Sec. 10).
- We introduce the selection of the strategy  $\hat{a}$  (see Sec. 11).
- We discuss the training efficiency of the proposed method (see Sec. 12).

## 1. Architecture of the Strategy Network

The architecture of the strategy network is illustrated in Fig. 1. We exploit ResNet18 [3] as the backbone. Given an image, the strategy network outputs an attack strategy, *i.e.*, the configuration of how to perform the adversarial attack. Specifically, the strategy network outputs a set of attack parameters for AE generation.

## 2. Proof of Theorem 1

First we introduce some notations. Let  $\mathcal{L}_0 : \mathcal{X} \times \mathcal{Y} \times \mathcal{W} \times \Theta \rightarrow \mathbb{R}^+$  be the objective function in (7) of submission (Line 400) as

$$\mathcal{L}_0 = \mathcal{L}_1 + \alpha \mathcal{L}_2 + \beta \mathcal{L}_3. \quad (1)$$

We define  $\mathbf{x}_{\text{adv}}^*(\mathbf{x}, \mathbf{w})$  as the optimal adversarial example generated by the strategy network

$$\begin{aligned} \mathbf{x}_{\text{adv}}^*(\mathbf{x}, \mathbf{w}) &= \arg \max_{\boldsymbol{\theta}} g(\mathbf{x}, \mathbf{a}(\boldsymbol{\theta}), \mathbf{w}) \\ &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{a} \sim p(\mathbf{a}|\mathbf{x}, \boldsymbol{\theta})} [\mathcal{L}_0], \end{aligned} \quad (2)$$

\*The first two authors contribute equally to this work. † Work done during an internship at Tencent AI Lab. ‡ Correspondence to: Baoyuan Wu (wubaoyuan@cuhk.edu.cn) and Xiaochun Cao (caoxiaochun@iie.ac.cn).

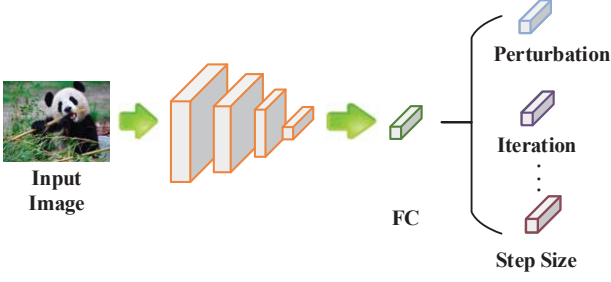


Figure 1. The overview of the strategy network. The outputs are a set of attack parameters for generating adversarial examples. Each attack parameter is encoded by a one-hot vector.

and  $\hat{x}_{\text{adv}}(\mathbf{x}, \mathbf{w})$  is a  $\delta$ -approximate solution to  $x_{\text{adv}}^*(\mathbf{x}, \mathbf{w})$ . In addition, the full gradient of  $\mathcal{L}_0$  w.r.t  $\mathbf{w}$  is

$$\begin{aligned}\nabla_{\mathbf{w}}\mathcal{L}_0(\mathbf{w}) &= \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}}\mathcal{L}_0^n \\ &= \frac{1}{N} \sum_{n=1}^N \nabla_{\mathbf{w}}\mathcal{L}_0(x_{\text{adv}}^*(\mathbf{x}_n, \mathbf{w}), \mathbf{w}),\end{aligned}\quad (3)$$

where  $x_{\text{adv}}^*(\mathbf{x}_n)$  is the optimal adversarial example for  $\mathbf{x}_n$ . The stochastic gradient of  $\mathcal{L}_0$  w.r.t  $\mathbf{w}$  is

$$\begin{aligned}\nabla_{\mathbf{w}}\ell(\mathbf{w}) &= \frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \nabla_{\mathbf{w}}\mathcal{L}_0^i \\ &= \frac{1}{|\mathcal{B}|} \sum_{n=1}^N \nabla_{\mathbf{w}}\mathcal{L}_0(x_{\text{adv}}^*(\mathbf{x}_i, \mathbf{w}), \mathbf{w}).\end{aligned}\quad (4)$$

Then  $\nabla_{\theta}\mathcal{L}_0$  and  $\nabla_{\theta}\ell$  correspond to the full and stochastic gradients of  $\mathcal{L}_0$  w.r.t  $\theta$ . Without loss of generality, we assume that

$$\mathbb{E}[\nabla_{\mathbf{w}}\ell(\mathbf{w})] = \nabla_{\mathbf{w}}\mathcal{L}_0(\mathbf{w}). \quad (5)$$

We note the approximate stochastic gradient as  $\nabla_{\mathbf{w}}\hat{\ell}$ :

$$\begin{aligned}\nabla_{\mathbf{w}}\hat{\ell}(\mathbf{w}) &= \frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \nabla_{\mathbf{w}}\hat{\mathcal{L}}_0^i \\ &= \frac{1}{|\mathcal{B}|} \sum_{n=1}^N \nabla_{\mathbf{w}}\mathcal{L}_0(\hat{x}_{\text{adv}}(\mathbf{x}_i, \mathbf{w}), \mathbf{w}).\end{aligned}\quad (6)$$

Moreover, the adversarial example  $x_{\text{adv}}(\mathbf{x}, \mathbf{w})$  can be identified by a parameter  $\theta$  of the strategy network and the gradients like (3), (4), (6) would be

$$\begin{aligned}\nabla_{\mathbf{w}}\mathcal{L}_0(\theta, \mathbf{w}) &:= \nabla_{\mathbf{w}}\mathcal{L}_0(\mathbf{w}) \\ \nabla_{\mathbf{w}}\ell(\theta, \mathbf{w}) &:= \nabla_{\mathbf{w}}\ell(\mathbf{w}) \\ \nabla_{\mathbf{w}}\hat{\ell}(\theta, \mathbf{w}) &:= \nabla_{\mathbf{w}}\hat{\ell}(\mathbf{w}).\end{aligned}\quad (7)$$

The corresponding gradients w.r.t  $\theta$  will be  $\nabla_{\theta}\mathcal{L}_0$ ,  $\nabla_{\theta}\ell$  and  $\nabla_{\theta}\hat{\ell}$ . As the  $\mathcal{L}_0$  in (7) of submission (Line 400) satisfies the

Lipschitz gradient conditions, given  $\mathbf{x}_n \in \mathcal{X}$ , it holds that

$$\begin{aligned}&\sup_{\theta} \|\nabla_{\mathbf{w}}\mathcal{L}_0^n(\theta, \mathbf{w}) - \nabla_{\mathbf{w}}\mathcal{L}_0^n(\theta, \mathbf{w}')\|_2 \\ &\leq L_{\mathbf{w}\mathbf{w}}\|\mathbf{w} - \mathbf{w}'\|_2 \\ &\sup_{\mathbf{w}} \|\nabla_{\mathbf{w}}\mathcal{L}_0^n(\theta, \mathbf{w}) - \nabla_{\mathbf{w}}\mathcal{L}_0^n(\theta', \mathbf{w})\|_2 \\ &\leq L_{\mathbf{w}\theta}\|\theta - \theta'\|_2 \\ &\sup_{\theta} \|\nabla_{\theta}\mathcal{L}_0^n(\theta, \mathbf{w}) - \nabla_{\theta}\mathcal{L}_0^n(\theta, \mathbf{w}')\|_2 \\ &\leq L_{\theta\mathbf{w}}\|\mathbf{w} - \mathbf{w}'\|_2,\end{aligned}\quad (8)$$

where  $L_{\mathbf{w}\mathbf{w}}$ ,  $L_{\mathbf{w}\theta}$  and  $L_{\theta\mathbf{w}}$  are positive constants. Furthermore, by the strongly-concavity of  $\mathcal{L}_0$  and given  $\mathbf{x}_n \in \mathcal{X}$ , we know that for any  $\theta_1$  and  $\theta_2 \in \Theta$ ,

$$\begin{aligned}&\mathcal{L}_0^n(\theta_1, \mathbf{w}) - \mathcal{L}_0^n(\theta_2, \mathbf{w}) \\ &\leq \langle \nabla_{\theta}\mathcal{L}_0^n(\theta, \mathbf{w}), \theta_1 - \theta_2 \rangle - \frac{\mu}{2}\|\theta_1 - \theta_2\|_2^2.\end{aligned}\quad (9)$$

As the variance of the stochastic gradient is bounded by  $\sigma^2 > 0$ , it means that

$$\mathbb{E}[\|\nabla_{\mathbf{w}}\ell(\mathbf{w}) - \nabla_{\mathbf{w}}\mathcal{L}_0(\mathbf{w})\|_2^2] \leq \sigma^2. \quad (10)$$

To prove the main result, we need the following two important lemmas.

**Lemma 1.** Suppose that  $\mathcal{L}_0$  in (7) of submission (Line 400) satisfies the Lipschitz gradient conditions as (8) and  $\mathcal{L}_0$  is  $\mu$ -strongly concave in  $\Theta$ , we have  $\mathcal{L}_0$  is Lipschitz smooth with  $L_0$

$$L_0 = \frac{L_{\mathbf{w}\theta}L_{\theta\mathbf{w}}}{\mu} + L_{\mathbf{w}\mathbf{w}}. \quad (11)$$

It holds that

$$\begin{aligned}\mathcal{L}_0(\mathbf{w}_1) &\leq \mathcal{L}_0(\mathbf{w}_2) + \langle \nabla_{\mathbf{w}}\mathcal{L}_0(\mathbf{w}_2), \mathbf{w}_1 - \mathbf{w}_2 \rangle \\ &\quad + \frac{L_0}{2}\|\mathbf{w}_1 - \mathbf{w}_2\|_2^2,\end{aligned}\quad (12)$$

and

$$\|\nabla_{\mathbf{w}}\mathcal{L}_0 - \nabla_{\mathbf{w}}\mathcal{L}_0(\mathbf{w}_2)\|_2 \leq L_0\|\mathbf{w}_1 - \mathbf{w}_2\|_2. \quad (13)$$

*Proof.* By the strongly-concavity of  $\mathcal{L}_0$  and given  $\mathbf{x}_n \in \mathcal{X}$ , for any  $\theta_1, \theta_2$  and the corresponding  $\mathbf{w}_1, \mathbf{w}_2$ , we have

$$\begin{aligned}&\mathcal{L}_0^n(\theta_1, \mathbf{w}_2) - \mathcal{L}_0^n(\theta_2, \mathbf{w}_2) \\ &\leq \langle \nabla_{\theta}\mathcal{L}_0^n(\theta_2, \mathbf{w}_2), \theta_1 - \theta_2 \rangle - \frac{\mu}{2}\|\theta_1 - \theta_2\|_2^2 \\ &\leq -\frac{\mu}{2}\|\theta_1 - \theta_2\|_2^2.\end{aligned}\quad (14)$$

The second inequality is true as

$$\langle \nabla_{\theta}\mathcal{L}_0^n(\theta_2, \mathbf{w}_2), \theta_1 - \theta_2 \rangle \leq 0.$$

In addition, we have

$$\begin{aligned} & \mathcal{L}_0^n(\theta_2, \mathbf{w}_2) - \mathcal{L}_0^n(\theta_1, \mathbf{w}_2) \\ & \leq \langle \nabla_{\theta} \mathcal{L}_0^n(\theta_1, \mathbf{w}_2), \theta_2 - \theta_1 \rangle - \frac{\mu}{2} \|\theta_1 - \theta_2\|_2^2 \quad (15) \\ & \leq -\frac{\mu}{2} \|\theta_1 - \theta_2\|_2^2. \end{aligned}$$

Combining (14) and (15), we have

$$\begin{aligned} & \mu \|\theta_1 - \theta_2\|_2^2 \\ & \leq \langle \nabla_{\theta} \mathcal{L}_0^n(\theta_1, \mathbf{w}_2), \theta_2 - \theta_1 \rangle \\ & \leq \langle \nabla_{\theta} \mathcal{L}_0^n(\theta_1, \mathbf{w}_2) - \nabla_{\theta} \mathcal{L}_0^n(\theta_1, \mathbf{w}_1), \theta_2 - \theta_1 \rangle \\ & \leq \|\nabla_{\theta} \mathcal{L}_0^n(\theta_1, \mathbf{w}_2) - \nabla_{\theta} \mathcal{L}_0^n(\theta_1, \mathbf{w}_1)\|_2 \|\theta_2 - \theta_1\|_2 \\ & \leq L_{\theta \mathbf{w}} \|\mathbf{w}_2 - \mathbf{w}_1\|_2 \|\theta_2 - \theta_1\|_2, \quad (16) \end{aligned}$$

where the second inequality holds as

$$\langle \nabla_{\theta} \mathcal{L}_0^n(\theta_1, \mathbf{w}_1), \theta_2 - \theta_1 \rangle \leq 0,$$

the third inequality follows from the Cauchy-Schwarz inequality, and the last one holds by the Lipschitz smoothness of the gradients of  $\mathcal{L}_0$  (8).

For any  $n \in [N]$ , we have

$$\begin{aligned} & \|\nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_1, \mathbf{w}_1) - \nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_2, \mathbf{w}_2)\|_2 \\ & \leq \|\nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_1, \mathbf{w}_1) - \nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_2, \mathbf{w}_1)\|_2 \\ & \quad + \|\nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_2, \mathbf{w}_1) - \nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_2, \mathbf{w}_2)\|_2 \quad (17) \\ & \leq L_{\mathbf{w}\theta} \|\theta_1 - \theta_2\|_2 + L_{\mathbf{w}\mathbf{w}} \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \\ & = \left( \frac{L_{\mathbf{w}\theta} L_{\theta \mathbf{w}}}{\mu} + L_{\mathbf{w}\mathbf{w}} \right) \|\mathbf{w}_1 - \mathbf{w}_2\|_2, \end{aligned}$$

where the first inequality follows from the triangle inequality, and the second inequality holds due to (16) and the Lipschitz smoothness of the gradients of  $\mathcal{L}_0$  (8). By the definition of  $\mathcal{L}$ , it holds that

$$\begin{aligned} & \|\nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}_1) - \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}_2)\|_2 \\ & = \left\| \frac{1}{N} \sum_{n=1}^N (\nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_1, \mathbf{w}_1) - \nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_2, \mathbf{w}_2)) \right\|_2 \\ & \leq \frac{1}{N} \sum_{n=1}^N \|\nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_1, \mathbf{w}_1) - \nabla_{\mathbf{w}} \mathcal{L}_0^n(\theta_2, \mathbf{w}_2)\|_2 \\ & \leq \left( \frac{L_{\mathbf{w}\theta} L_{\theta \mathbf{w}}}{\mu} + L_{\mathbf{w}\mathbf{w}} \right) \|\mathbf{w}_1 - \mathbf{w}_2\|_2. \quad (18) \end{aligned}$$

With the definition of the Lipschitz smoothness, we complete the proof.  $\square$

**Lemma 2.** Suppose that  $\mathcal{L}_0$  in (7) of submission (Line 400) satisfies the Lipschitz gradient conditions as (8) and  $\mathcal{L}_0$  is  $\mu$ -strongly concave in  $\Theta$ , the approximate stochastic gradient  $\nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w})$  (6) satisfies

$$\|\nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}) - \nabla_{\mathbf{w}} \ell(\mathbf{w})\|_2 \leq L_{\mathbf{w}\theta} \sqrt{\frac{\delta}{\mu}}, \quad (19)$$

where  $\hat{\mathbf{x}}_{\text{adv}}(\mathbf{x}, \mathbf{w})$  is a  $\delta$ -approximate solution to  $\mathbf{x}_{\text{adv}}^*(\mathbf{x}, \mathbf{w})$  with given  $\mathbf{x} \in \mathcal{X}$ .

*Proof.* By the definitions of  $\nabla_{\mathbf{w}} \hat{\ell}$  and  $\nabla_{\mathbf{w}} \ell$ , we have

$$\begin{aligned} & \|\nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}) - \nabla_{\mathbf{w}} \ell(\mathbf{w})\|_2 \\ & = \left\| \frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} (\nabla_{\mathbf{w}} \mathcal{L}_0(\hat{\mathbf{x}}_{\text{adv}}(\mathbf{x}_i, \mathbf{w}), \mathbf{w}) - \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{x}_{\text{adv}}^*(\mathbf{x}_i, \mathbf{w}), \mathbf{w})) \right\|_2 \\ & \leq \frac{1}{|\mathcal{B}|} \sum_{n=1}^N \|\nabla_{\mathbf{w}} \mathcal{L}_0(\hat{\mathbf{x}}_{\text{adv}}(\mathbf{x}_i, \mathbf{w}), \mathbf{w}) - \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{x}_{\text{adv}}^*(\mathbf{x}_i, \mathbf{w}), \mathbf{w})\|_2 \\ & \leq \frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} L_{\mathbf{w}\theta} \|\hat{\theta} - \theta^*\|_2, \quad (20) \end{aligned}$$

where the second inequality follows from the triangle inequality, the third inequality holds due to the gradient Lipschitz condition, and  $\hat{\theta}$  is the parameter of strategy network corresponding to  $\hat{\mathbf{x}}_{\text{adv}}(\mathbf{x}_i, \mathbf{w})$ ,  $\theta^*$  is similar.

Since  $\hat{\mathbf{x}}_{\text{adv}}(\mathbf{x}_i, \mathbf{w})$  is a  $\delta$ -approximate adversarial example generated by the strategy network, we have

$$\langle \theta^* - \hat{\theta}, \nabla_{\theta} \mathcal{L}_0(\hat{\theta}, \mathbf{w}) \rangle \leq \delta. \quad (21)$$

In addition, it holds that

$$\langle \hat{\theta} - \theta^*, \nabla_{\theta} \mathcal{L}_0(\theta^*, \mathbf{w}) \rangle \leq 0. \quad (22)$$

Putting (21) and (22) together gives birth to

$$\langle \hat{\theta} - \theta^*, \nabla_{\theta} \mathcal{L}_0(\theta^*, \mathbf{w}) - \nabla_{\theta} \mathcal{L}_0(\hat{\theta}, \mathbf{w}) \rangle \leq \delta. \quad (23)$$

Moreover, by the strongly concavity of  $\mathcal{L}_0$  and (16), we have

$$\begin{aligned} & \mu \|\theta^* - \hat{\theta}\|_2^2 \\ & \leq \langle \nabla_{\theta} \mathcal{L}_0^n(\theta^*, \mathbf{w}) - \nabla_{\theta} \mathcal{L}_0^n(\hat{\theta}, \mathbf{w}), \hat{\theta} - \theta^* \rangle \quad (24) \\ & \leq \delta. \end{aligned}$$

Consequently, it immediately yields

$$\|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}\|_2 \leq \sqrt{\frac{\delta}{\mu}}. \quad (25)$$

Substituting (25) into (20), we complete the proof.  $\square$

*Proof.* By Lemma 1, we have

$$\begin{aligned} \mathcal{L}_0(\mathbf{w}^{t+1}) &\leq \mathcal{L}_0(\mathbf{w}^t) + \frac{L_0}{2} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|_2^2 \\ &\quad + \langle \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t), \mathbf{w}^{t+1} - \mathbf{w}^t \rangle. \end{aligned}$$

Due to

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t),$$

it holds that

$$\begin{aligned} &\mathcal{L}_0(\mathbf{w}^{t+1}) \\ &\leq \mathcal{L}_0(\mathbf{w}^t) - \eta_t \|\nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 + \frac{L_0 \eta_t^2}{2} \|\nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t)\|_2^2 \\ &\quad + \eta_t \langle \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t), \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t) - \nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t) \rangle \\ &= \mathcal{L}_0(\mathbf{w}^t) - \eta_t \left(1 - \frac{L_0 \eta_t}{2}\right) \|\nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 \\ &\quad + \eta_t \left(1 - \frac{L_0 \eta_t}{2}\right) \langle \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t), \\ &\quad \quad \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t) - \nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t) \rangle \\ &\quad + \frac{L_0 \eta_t^2}{2} \|\nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t) - \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 \\ &= \mathcal{L}_0(\mathbf{w}^t) - \eta_t \left(1 - \frac{L_0 \eta_t}{2}\right) \|\nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 \\ &\quad + \eta_t \left(1 - \frac{L_0 \eta_t}{2}\right) \langle \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t), \\ &\quad \quad \nabla_{\mathbf{w}} \ell(\mathbf{w}^t) - \nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t) \rangle \\ &\quad + \eta_t \left(1 - \frac{L_0 \eta_t}{2}\right) \langle \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t), \\ &\quad \quad \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t) - \nabla_{\mathbf{w}} \ell(\mathbf{w}^t) \rangle \\ &\quad + \frac{L_0 \eta_t^2}{2} \|\nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t) - \nabla_{\mathbf{w}} \ell(\mathbf{w}^t) \\ &\quad \quad + \nabla_{\mathbf{w}} \ell(\mathbf{w}^t) - \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 \\ &\leq \mathcal{L}_0(\mathbf{w}^t) - \frac{\eta_t}{2} \left(1 - \frac{L_0 \eta_t}{2}\right) \|\nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 \\ &\quad + \frac{\eta_t}{2} \left(1 - \frac{L_0 \eta_t}{2}\right) \|\nabla_{\mathbf{w}} \ell(\mathbf{w}^t) - \nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t)\|_2^2 \\ &\quad + \eta_t \left(1 + \frac{L_0 \eta_t}{2}\right) \langle \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t), \\ &\quad \quad \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t) - \nabla_{\mathbf{w}} \ell(\mathbf{w}^t) \rangle \\ &\quad + L_0 \eta_t^2 \|\nabla_{\mathbf{w}} \hat{\ell}(\mathbf{w}^t) - \nabla_{\mathbf{w}} \ell(\mathbf{w}^t)\|_2^2 \\ &\quad + L_0 \eta_t^2 \|\nabla_{\mathbf{w}} \ell(\mathbf{w}^t) - \nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 \end{aligned} \quad (26)$$

Taking expectation on both sides of the above inequality

conditioned on  $\mathbf{w}^t$ , then we have

$$\begin{aligned} & \mathbb{E}[\mathcal{L}_0(\mathbf{w}^{t+1}) - \mathcal{L}_0(\mathbf{w}^t) | \mathbf{w}^t] \\ & \leq -\frac{\eta_t}{2} \left(1 - \frac{L_0 \eta_t}{2}\right) \|\nabla_{\mathbf{w}} \mathcal{L}_0(\mathbf{w}^t)\|_2^2 \\ & \quad + \frac{\eta_t}{2} \left(1 + \frac{3\eta_t L_0}{2}\right) \frac{\delta L_{\mathbf{w}\theta}^2}{\mu} + L_0 \eta_t^2 \sigma^2. \end{aligned} \quad (27)$$

Then we do the telescope sum over  $t = 0, \dots, T-1$ , we obtain

$$\begin{aligned} & \sum_{t=0}^{T-1} \frac{\eta_t}{2} \left(1 - \frac{L_0 \eta_t}{2}\right) \mathbb{E}[\|\mathcal{L}_0(\mathbf{w}^t)\|_2^2] \\ & \leq \mathbb{E}[\mathcal{L}_0(\mathbf{w}^0) - \mathcal{L}_0(\mathbf{w}^T)] + L_0 \sum_{t=0}^{T-1} \eta_t^2 \sigma^2 \\ & \quad + \sum_{t=0}^{T-1} \frac{\eta_t}{2} \left(1 + \frac{3\eta_t L_0}{2}\right) \frac{\delta L_{\mathbf{w}\theta}^2}{\mu}. \end{aligned} \quad (28)$$

Choosing  $\eta_t = \eta_1$  as

$$\eta_1 = \min \left( \frac{1}{L_0}, \sqrt{\frac{\mathcal{L}_0(\mathbf{w}^0) - \min_{\mathbf{w}} \mathcal{L}_0(\mathbf{w})}{\sigma^2 T L_0}} \right), \quad (29)$$

it holds that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla \mathcal{L}_0(\mathbf{w}^t)\|_2^2] \leq 4\sigma \sqrt{\frac{\Delta L_0}{T}} + \frac{5\delta L_{\mathbf{w}\theta}^2}{\mu}, \quad (30)$$

where  $\Delta = \mathcal{L}_0(\mathbf{w}^0) - \min_{\mathbf{w}} \mathcal{L}_0(\mathbf{w})$ .  $\square$

### 3. Details of the databases

CIFAR-10, CIFAR-100, and Tiny ImageNet are the most widely used databases for the evaluation of adversarial robustness. The CIFAR-10 dataset contains 50,000 training images and 10,000 test images, which covers 10 classes of images in the size of  $32 \times 32$ . The CIFAR-100 dataset also contains 50,000 training images and 10,000 test images in the size of  $32 \times 32$ , but it covers 100 classes. The Tiny ImageNet database is a subset collected from the ImageNet database, which covers 200 classes with 600 images in the size of  $64 \times 64$  for each class. As there are no labels for test images of Tiny ImageNet, following [4], we evaluate methods on the validation set.

### 4. Detailed Settings for LAS-AT

**LAS-PGD-AT:** The proposed LAS-PGD-AT is implemented based on PGD-AT [6], which performs the early

stopping adversarial training. PGD-AT is trained on the adversarial examples generated by a fixed PGD attack strategy. It has three stationary attack parameters that depict how to attack, *i.e.*, the maximal perturbation strength, the attack step, and the attack iteration. Hence, the strategy network of our LAS-PGD-AT has three parallel softmax layers to predict the three attack parameters. Following the default setting in [6], we adopt SGD with momentum 0.9, weight decay  $5 \times 10^{-4}$ , and batch size 128. LAS-PGD-AT and PGD-AT [6] are trained for 110 epochs. The learning rate decays with a factor of 0.1 at the 100 and 105 epochs, respectively.

**LAS-TRADES:** The proposed LAS-TRADES is implemented based on TRADES [10] which exploits a regularized surrogate loss to perform adversarial training. Since TRADES uses the same three attack parameters as PGD-AT to depict how to attack, the strategy network of LAS-TRADES shares the same output configuration as LAS-PGD-AT. Different from LAS-PGD-AT, LAS-TRADES is trained with the regularized surrogate loss rather than the cross-entropy loss. Following the default setting in [10], we adopt SGD with momentum 0.9, weight decay  $2 \times 10^{-4}$ , and batch size 128. LAS-TRADES and TRADES [10] are trained for 100 epochs. The learning rate decays with a factor of 0.1 at the 75 and 90 epochs, respectively.

**LAS-AWP:** The proposed LAS-AWP is implemented based on AWP [9]. AWP has 3 stationary attack parameters, *i.e.*, the maximal perturbation strength, the attack step, and the attack iteration. Hence, the strategy network of LAS-AWP has 3 parallel softmax layers to predict these attack parameters. Following the default setting in [9], we adopt SGD with momentum 0.9, weight decay  $5 \times 10^{-4}$ , and batch size 128. LAS-AWP and AWP [9] are trained for 200 epochs. The learning rate decays with a factor of 0.1 at the 100 and 150 epochs, respectively.

### 5. Detailed Training And Evaluation Settings

In Table 7 of the manuscript, we compared our LAS-Madry-AT with CAT [1], DART [7], and FAT [11]. We adopt the default setting in their original paper to train them. As for LAS-Madry-AT, it is trained for 200 epochs. The learning rate decays with a factor of 0.1 at the 100 and 150 epochs, respectively. We report the performance at the last epoch. In Fig. 3 of the manuscript, we compared our LAS-PGD-AT with OHL [5] and AdvHP [12]. OHL [5] and AdvHP [12] search the hyper-parameter of data augmentation for image classification. For a fair comparison, we use the same hyper-parameters and search range for them as ours. Specifically, the range of the maximal perturbation strength is set from 3 to 15, the range of the attack step is set from 1 to 6, and the range of the attack iteration is set from 3 to 15. We also adopt the same strategy network and training setting for them. Specifically, for the target network, we adopt

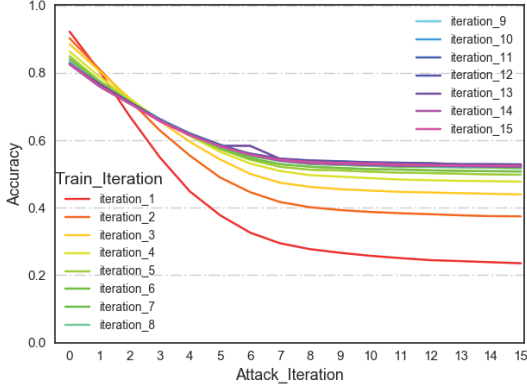


Figure 2. The influence of the number of iteration in PGD-AT [6]. Curves represent performance of PGD-AT [6] with different numbers of iteration for training. X-axis represents the number of iteration of the PGD attack. Y-axis represents the accuracy of adversarial examples generated by the PGD attack method.

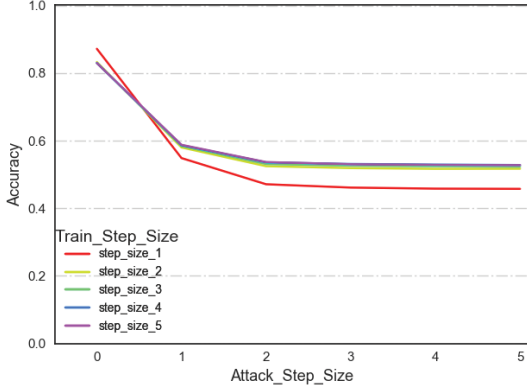


Figure 3. The influence of the step size in PGD-AT [6]. Curves represent performance of PGD-AT [6] with different step sizes for training. X-axis represents the step size of the PGD attack. Y-axis represents the accuracy of adversarial examples generated by the PGD attack method.

SGD momentum optimizer with a learning rate of 0.1, and a weight decay of  $5 \times 10^{-4}$ . ResNet18 is used as the backbone of the strategy network. The training epoch is set to 110. The learning rate decays with a factor of 0.1 at the 105 and 110 epochs, respectively. In this way, the difference between our proposed method and them is the loss term which guides the learning of the strategy network.

## 6. More Comparisons with AWP

**Comparisons with AWP trained with more iterations and a larger perturbation strength.** To evaluate the proposed method’s effectiveness in improving the model robustness, we compare our LAS-AWP with the AWP [9] trained with more iterations ( $I_{\text{train}} = 15$ ) and a larger per-

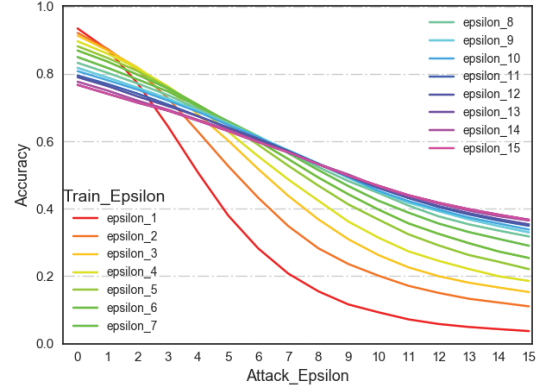


Figure 4. The influence of the maximal perturbation strength in PGD-AT [6]. Curves represent performance of PGD-AT [6] with different maximal perturbation strengths for training. X-axis represents the perturbation strength of the PGD attack. Y-axis represents the accuracy of adversarial examples generated by the PGD attack method.

turbation strength ( $\epsilon_{\text{train}} = 15$ ). The results are shown in Table 1. Our LAS-AWP not only achieves a higher robust accuracy under all attack scenarios but also achieves a higher clean accuracy on clean images. We can attribute the improvements to using automatically generated attack strategies instead of more iterations and a larger perturbation strength.

## 7. Influence of Attack Parameters

In order to explore the influence of the maximal perturbation strength, the number of iteration, and the step size in adversarial training, we conduct an experiment on CIFAR10 by using the PGD-AT [6] with different maximal perturbation strengths, different numbers of iteration, or different step sizes to perform adversarial training. During testing, we use PGD attack with different perturbation strengths, different numbers of iteration, or different step sizes to evaluate the robustness of the trained target model. We use ResNet18 as the target model.

**To explore the influence of the maximal perturbation strength,** we use the PGD-AT [6] with different maximal perturbation strengths to train the target model. We then use PGD attack with different perturbation strengths to attack the target model. The number of iteration of the PGD attack is set to 10 and the step size is set to 2. Note that most works use the PGD attack method with the perturbation strength of 8, iteration of 10, and step size of 2 (PGD(8,10,2)) to evaluate model robustness. The results are shown in Fig. 4. ‘Train\_Epsilon’ represents the maximal perturbation strength used in adversarial training while ‘Attack\_Epsilon’ represents the perturbation strength of the PGD attack for evaluation.



Table 1. Test robustness (%) on the CIFAR-10 database using ResNet18. Number in bold indicates the best.

Method	Clean	PGD-10	PGD-20	PGD-50	C&W	AA
AWP( $I_{\text{train}} = 10, \epsilon_{\text{train}} = 8$ )	80.72	55.33	54.78	54.28	51.67	49.44
AWP( $I_{\text{train}} = 10, \epsilon_{\text{train}} = 15$ )	66.73	52.24	52.14	52.06	48.1	47.03
AWP( $I_{\text{train}} = 15, \epsilon_{\text{train}} = 8$ )	80.13	55.82	55.24	55.13	51.53	49.62
LAS-AWP(ours)	<b>83.03</b>	<b>56.45</b>	<b>55.76</b>	<b>55.43</b>	<b>53.06</b>	<b>50.77</b>

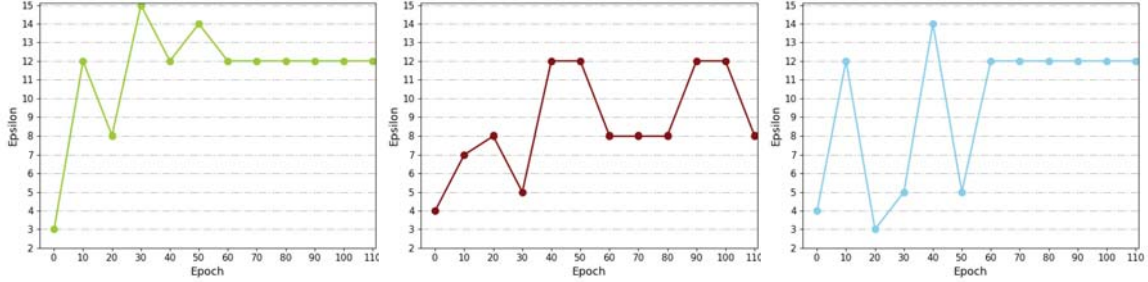


Figure 5. The evolution of the generated perturbation strength of several samples during the whole training process. X-axis represents the training epoch. Y-axis represents the perturbation strength.

Table 2. Hyperparameter selection. Test robustness (%) on the CIFAR-10 database using ResNet18. Number in bold indicates the best.

Method	Clean	PGD-10	AA
$\alpha = 2$	$\beta = 2$	<b>82.32</b>	54.27
	$\beta = 4$	82.05	<b>54.29</b>
	$\beta = 8$	82.01	54.28
$\alpha = 4$	$\beta = 2$	82.02	54.36
	$\beta = 4$	81.64	54.05
	$\beta = 8$	81.92	54.19
$\alpha = 8$	$\beta = 2$	81.73	54.10
	$\beta = 4$	82.18	54.30
	$\beta = 8$	82.10	54.35

Observations are summarized as follows. **First**, though using a large maximal perturbation strength for training improves the robustness, it decreases the clean accuracy. For example, when  $\text{Train\_Epsilon} = 15$ , the clean accuracy (*i.e.*,  $\text{Attack\_Epsilon} = 0$ ) drops to 76.7% which is much worse than that of  $\text{Train\_Epsilon} \leq 8$ . **Second**, using a large  $\text{Train\_Epsilon}$  can achieve better performance than using a small one only when the  $\text{Attack\_Epsilon}$  is large. Comparing the performance of  $\text{Train\_Epsilon} = 8$  and  $\text{Train\_Epsilon}$

$= 15$ , the accuracy of  $\text{Train\_Epsilon} = 15$  is lower when  $\text{Attack\_Epsilon} < 8$ , while the accuracy is higher when  $\text{Attack\_Epsilon} > 8$ . **Third**, when  $\text{Attack\_Epsilon} = 8$ , increasing  $\text{Train\_Epsilon}$  could slightly improve the robustness, but gets sharp clean accuracy drop.

As mentioned above, when the perturbation strength of PGD attack is 8 for evaluation (*i.e.*,  $\text{Attack\_Epsilon} = 8$ ), using a large maximal perturbation strength can only slightly improve the robustness, but will hurt the clean accuracy a lot. **Hence, in the manuscript, the maximal perturbation strengths of other state-of-the-art adversarial training methods are set to the same as their original papers, *i.e.*,  $\text{Train\_Epsilon} = 8$ .** The range of the maximal perturbation strength is set to  $[3, 15]$  in our method.

**To explore the influence of the number of iteration**, we use the PGD-AT [6] with different numbers of iteration to train the target model. We then use PGD attack with different numbers of iteration to attack the target model. The perturbation strength of the PGD attack is set to 8 and the step size is set to 2. The results are shown in Fig. 2. ‘ $\text{Train\_Iteration}$ ’ represents the number of iteration used in adversarial training while ‘ $\text{Attack\_Iteration}$ ’ represents the number of iteration of the PGD attack for evaluation. As shown in Fig. 2, we can observe a similar phenomenon as in Fig. 4. Increasing the number of iteration could improve the robustness, but it will hurt the clean accuracy. When  $\text{Train\_Iteration}$  is larger than 10, we can only slightly improve the robustness by increas-

ing Train.Iteration, but gets a slight drop in clean accuracy.

To explore the influence of the step size, we use the PGD-AT [6] with different step sizes to train the target model. We then use PGD attack with different step sizes to attack the target model. The maximal perturbation strength of the PGD attack is set to 8 and the number of iteration is set to 10. The results are shown in Fig. 3. ‘Train\_Step\_Size’ represents the step size used in adversarial training while ‘Attack\_Step\_Size’ represents the step size of the PGD attack for evaluation. As shown in Fig. 3, when Train\_Step\_Size is larger than 2, the clean accuracy is nearly the same, while the robustness has marginal changes.

## 8. Selection of Hyper-parameters

There are two trade-off hyper-parameters in the objective function of the proposed method, *i.e.*,  $\alpha$  and  $\beta$  in Eq. (8) in the manuscript. In this section, we present the performance of our proposed LAS-PGD-AT [6] with different  $\alpha$  and  $\beta$  pairs on CIFAR-10. The range of  $\alpha$  and  $\beta$  is set to  $\{2^1, 2^2, 2^3\}$ . The results are shown in Table 2.

When  $\alpha = 2$  and  $\beta = 2$ , LAS-PGD-AT [6] achieves the best performance in clean accuracy and PGD-10 attack. And when  $\alpha = 2$  and  $\beta = 4$ , the proposed method achieves the best performance in APGD-T, FAB, SQUARE, and AA attacks. And it also achieves the competitive performance in clean accuracy and PGD-10 attack. Moreover, our method is not sensitive to the two hyper-parameters as the robustness and the clean accuracy does not change in a large range. Hence, in this paper, we set  $\alpha$  to 2 and  $\beta$  to 4.

## 9. Illustration of the Evolution of Generated Perturbation Strength

Given a sample, our strategy network generates an attack strategy. As the maximal perturbation strength affects the performance the most, we illustrate how the generated perturbation strength changes during adversarial training. Fig. 5 presents the evolution of the generated perturbation strength of three randomly selected images during the whole training process.

It can be observed that the perturbation strength of the same sample changes dynamically during the training process. At the beginning of adversarial training, the perturbation strength is small. As the training process goes on and the robustness of the target network improves, the perturbation strength becomes larger.

## 10. Experiments on more databases

We conduct experiments under the same  $L_\infty = 8$  as the manuscript. 1) For GTSRB, following [2], we adopt

Table 3. Test robustness (%) on GTSRB using ResNet18.

Method	Clean	PGD-50	C&W	AA
Clean	<b>98.36</b>	15.60	16.31	13.07
PGD-AT [6]	90.49	62.11	62.83	60.80
TRADES [10]	88.17	63.02	62.86	61.35
AWP [9]	92.68	63.85	64.65	61.66
LAS-AT(ours)	92.27	64.98	64.48	62.64
LAS-TRADES(ours)	90.26	64.62	64.14	62.59
LAS-AWP(ours)	93.79	<b>66.68</b>	<b>67.56</b>	<b>64.44</b>

Table 4. Test robustness (%) on ImageNet using ResNet50.

Method	Clean	PGD-50	C&W	AA
PGD-AT	47.86	23.63	23.04	15.23
LAT-AT(ours)	<b>48.03</b>	<b>25.97</b>	<b>24.21</b>	<b>16.56</b>

Table 5. Test robustness (%) on CIFAR-10 using WRN34-10.

$\hat{a}$	Clean	PGD-50	C&W	AA
PGD-10	86.23	56.12	55.73	53.58
PGD-20	86.19	56.14	<b>55.76</b>	53.68
PGD-50	86.2	<b>56.21</b>	55.68	<b>53.75</b>
C&W	<b>86.29</b>	55.94	55.48	53.62

ResNet18 as backbone. As shown in Table 3, LAS-AT improves the robust accuracy of three base models under all attack scenarios and also improves their clean accuracy. 1) For ImageNet, following [8], we compare with PGD-AT using ResNet50 as backbone. As shown in Table 4, LAS-AT achieves the better clean and robust accuracy.

## 11. Selection of the Strategy

The strategy  $\hat{a}$  is used to evaluate the robustness of the model updated with the current strategy, which can be any attack strategy. In our experiments,  $\hat{a}$  is set to the widely used attack strategy, *i.e.*, PGD-10. Here we conduct experiments with different attack strategies  $\hat{a}$  on CIFAR-10 with WRN34-10 as backbone, and find that our method is not affected too much by  $\hat{a}$ , as shown in Table. 5.

## 12. Discussion about the Training Efficiency

Compared to ordinary training, the most time-consuming part of AT is adversarial example (AE) generation. 1) As claimed in Sec. 3.4, our method calls AE generation to **update the target network** (T-net) with from 3 to 15 attack iterations (Line 588), and to **update the strategy network** (S-net) with 10 attack iterations. The S-net is updated once after  $k$ -times updates of the T-net. Thus, for each update of T-net, our LAS-AT takes  $15 + 10/k$  attack iterations at most, while PGD-AT takes 10 iterations. 2) In experiments, we set  $k = 40$  (Line 640). Thus, in the worst case, LAS-AT takes 52.5% more than PGD-AT. As shown in Tab. 1



of the manuscript, the actual additional training cost of our method is about 40% of PGD-AT’s cost. Our method can be scalable to large datasets.

## References

- [1] Qi-Zhi Cai, Min Du, Chang Liu, and Dawn Song. Curriculum adversarial training. *arXiv preprint arXiv:1805.04807*, 2018. 5
- [2] Jiefeng Chen, Xi Wu, Yang Guo, Yingyu Liang, and Somesh Jha. Towards evaluating the robustness of neural networks learned by transduction. *arXiv preprint arXiv:2110.14735*, 2021. 8
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1
- [4] Saehyung Lee, Hyungyu Lee, and Sungroh Yoon. Adversarial vertex mixup: Toward better adversarially robust generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 272–281, 2020. 5
- [5] Chen Lin, Minghao Guo, Chuming Li, Xin Yuan, Wei Wu, Junjie Yan, Dahua Lin, and Wanli Ouyang. Online hyper-parameter learning for auto-augmentation strategy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6579–6588, 2019. 5
- [6] Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In *International Conference on Machine Learning*, pages 8093–8104. PMLR, 2020. 1, 5, 6, 7, 8
- [7] Yisen Wang, Xingjun Ma, James Bailey, Jinfeng Yi, Bowen Zhou, and Quanzhan Gu. On the convergence and robustness of adversarial training. In *ICML*, volume 1, page 2, 2019. 5
- [8] Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training. *arXiv preprint arXiv:2001.03994*, 2020. 8
- [9] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. *Advances in Neural Information Processing Systems*, 33, 2020. 5, 6, 8
- [10] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. *arXiv preprint arXiv:1901.08573*, 2019. 5, 8
- [11] Jingfeng Zhang, Xilie Xu, Bo Han, Gang Niu, Lizhen Cui, Masashi Sugiyama, and Mohan Kankanhalli. Attacks which do not kill training make adversarial learning stronger. In *International Conference on Machine Learning*, pages 11278–11287. PMLR, 2020. 5
- [12] Xinyu Zhang, Qiang Wang, Jian Zhang, and Zhao Zhong. Adversarial autoaugment. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 5