

---

# Supplementary Material for: Adversarial Distributional Training for Robust Deep Learning

---

Yinpeng Dong\*, Zhijie Deng\*, Tianyu Pang, Jun Zhu, Hang Su†  
Dept. of Comp. Sci. & Tech., Institute for AI, BNRist Center  
Tsinghua-Bosch Joint ML Center, THBI Lab, Tsinghua University, Beijing, 100084 China  
{dyp17, dzj17, pty17}@mails.tsinghua.edu.cn, {suhangss, dcszj}@mail.tsinghua.edu.cn

## A Technical details and algorithms

### A.1 ADT<sub>EXP</sub>

We provide the algorithm for ADT<sub>EXP</sub> in Alg. 2.

### A.2 ADT<sub>EXP-AM</sub>

By amortizing the explicit adversarial distributions, we can rewrite the minimax problem of ADT as

$$\min_{\theta} \max_{\phi} \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{E}_{p_{\phi}(\delta_i|\mathbf{x}_i)} [\mathcal{L}(f_{\theta}(\mathbf{x}_i + \delta_i), y_i)] + \lambda \mathcal{H}(p_{\phi}(\delta_i|\mathbf{x}_i)) \right\}, \quad (\text{A.1})$$

where  $\theta$  and  $\phi$  are the parameters of the DNN classifier and the generator, respectively. During training, we perform stochastic gradient descent and ascent on  $\theta$  and  $\phi$  simultaneously, to accomplish adversarial training. To enable the gradients flowing from  $\delta_i$  to  $\phi$ , we apply the same reparameterization strategy as in Sec. 3.1. In practice, we only use one MC sample for each data. We provide the algorithm for ADT<sub>EXP-AM</sub> in Alg. 3.

### A.3 ADT<sub>IMP-AM</sub>

For the implicit adversarial distributions, we have no access to the density  $p_{\phi}(\delta_i|\mathbf{x}_i)$ , such that the entropy of the adversarial distributions cannot be estimated exactly<sup>1</sup>. An appealing alternative is to maximize the variational lower bound of the entropy [2] for its simplicity and success in GANs [3]. In our case, for a natural input  $\mathbf{x}_i$ , we can similarly derive the following lower bound stemming from the mutual information between the perturbation  $\delta_i$  and the random noise  $\mathbf{z}$  (proof in Appendix B.3) as

$$\mathcal{H}(p_{\phi}(\delta_i|\mathbf{x}_i)) \geq \mathcal{U}(q) = \mathbb{E}_{p(\mathbf{z})} \log q(\mathbf{z}|g_{\phi}(\mathbf{z}; \mathbf{x}_i)) + c, \quad (\text{A.2})$$

where  $c$  is a constant and  $q(\cdot|\cdot)$  is an introduced variational distribution. In practice, we implement  $q$  as a diagonal Gaussian, whose mean and standard derivation are given by a  $\psi$ -parameterized neural network. Then we have the training objective as

$$\min_{\theta} \max_{\phi, \psi} \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{E}_{p(\mathbf{z})} [\mathcal{L}(f_{\theta}(\mathbf{x}_i + g_{\phi}(\mathbf{z}; \mathbf{x}_i)), y_i) + \lambda \log q_{\psi}(\mathbf{z}|g_{\phi}(\mathbf{z}; \mathbf{x}_i))] \right\}, \quad (\text{A.3})$$

which is solved by simultaneous stochastic gradient descent and ascent on  $\theta$  and  $(\phi, \psi)$ . We provide the algorithm for ADT<sub>IMP-AM</sub> in Alg. 4.

---

\*Equal contribution. † Corresponding author.

<sup>1</sup>We can also directly estimate the gradient of the entropy with advanced techniques such as spectral Stein gradient estimator [15], and we leave this for future work.

---

**Algorithm 2** The training algorithm for ADT<sub>EXP</sub>

---

**Input:** Training data  $\mathcal{D}$ , objective function  $\mathcal{J}(p_{\phi_i}(\delta_i), \theta)$ , training epochs  $N$ , the number of inner maximization steps  $T$ , the number of MC samples for gradient estimation in each step  $k$ , and learning rates  $\eta_\theta, \eta_\phi$ .

```
1: Initialize  $\theta$ ;  
2: for epoch = 1 to  $N$  do  
3:   for each minibatch  $\mathcal{B} \subset \mathcal{D}$  do  
4:     for each input  $(\mathbf{x}_i, y_i) \in \mathcal{B}$  do  
5:       Initialize  $\phi_i$ ;  
6:       for  $t = 1$  to  $T$  do  
7:         Calculate the gradient  $\mathbf{g}_i$  of  $\phi_i$  by Eq. (10) via MC integration using  $k$  samples;  
8:         Update  $\phi_i$  with gradient ascent
```

$$\phi_i \leftarrow \phi_i + \eta_\phi \cdot \mathbf{g}_i.$$

```
9:     end for  
10:   end for  
11:   Update  $\theta$  with stochastic gradient descent
```

$$\theta \leftarrow \theta - \eta_\theta \cdot \mathbb{E}_{(\mathbf{x}_i, y_i) \in \mathcal{B}} [\nabla_\theta \mathcal{J}(p_{\phi_i}(\delta_i), \theta)].$$

```
12: end for  
13: end for
```

---

**Algorithm 3** The training algorithm for ADT<sub>EXP-AM</sub>

---

**Input:** Training data  $\mathcal{D}$ , objective function in Eq. (A.1), training epochs  $N$ , and learning rates  $\eta_\theta, \eta_\phi$ .

```
1: Initialize  $\theta$  and  $\phi$ ;  
2: for epoch = 1 to  $N$  do  
3:   for each minibatch  $\mathcal{B} \subset \mathcal{D}$  do  
4:     Input  $\mathbf{x}_i$  to the generator and obtain the distribution parameters  $(\mu_i, \sigma_i)$  for each data  
      $(\mathbf{x}_i, y_i) \in \mathcal{B}$ ;  
5:     Sample one  $\delta_i$  from the distribution defined by Eq. (8) given  $(\mu_i, \sigma_i)$  for each  $(\mathbf{x}_i, y_i) \in$   
      $\mathcal{B}$  to approximately calculate the gradient of Eq. (A.1) w.r.t.  $\theta$  and  $\phi$ , and obtain  $\mathbf{g}_\theta$  and  $\mathbf{g}_\phi$ ;  
6:     Update  $\theta$  by:  $\theta \leftarrow \theta - \eta_\theta \cdot \mathbf{g}_\theta$ .  
7:     Update  $\phi$  by:  $\phi \leftarrow \phi + \eta_\phi \cdot \mathbf{g}_\phi$ .  
8:   end for  
9: end for
```

---

## B Proofs

We provide the proofs in this section.

### B.1 Proof of Theorem 1

**Assumption 1.** The loss function  $\mathcal{J}(p(\delta_i), \theta)$  is continuously differentiable w.r.t.  $\theta$ .

**Assumption 2.** Probability density functions of distributions in  $\mathcal{P}$  are bounded and equicontinuous.

**Remark 1.** For the explicit adversarial distributions defined in Eq. (8), we can assume that the mean and standard deviation of each dimension satisfy  $|\mu_i^{(j)}| < \kappa_\mu$  and  $\kappa_\sigma^{lo} < \sigma_i^{(j)} < \kappa_\sigma^{up}$ , where  $\kappa_\mu, \kappa_\sigma^{lo}$ , and  $\kappa_\sigma^{up}$  are constants. Note that they can be easily satisfied since we add an entropic regularization term into the training objective (5), such that the mean cannot be too large while the standard deviation cannot be too small or too large given Eq. (11). In practice, we can clip  $\mu_i^{(j)}$  and  $\sigma_i^{(j)}$  if they are out of the thresholds. Then we can prove that the density functions of the explicit adversarial distributions defined in Eq. (8) are bounded and equicontinuous, satisfying Assumption 2. However, for the implicit adversarial distributions introduced in Sec. 3.3, we cannot prove that Assumption 2 is satisfied. Though unsatisfied, the experiments suggest that we can still rely on Theorem 1 and the general algorithm for training.

---

**Algorithm 4** The training algorithm for ADT<sub>IMP-AM</sub>


---

**Input:** Training data  $\mathcal{D}$ , objective function in Eq. (A.3), training epochs  $N$ , and learning rates  $\eta_\theta, \eta_\phi, \eta_\psi$ .

- 1: Initialize  $\theta, \phi$ , and  $\psi$ ;
  - 2: **for** epoch = 1 **to**  $N$  **do**
  - 3:     **for** each minibatch  $\mathcal{B} \subset \mathcal{D}$  **do**
  - 4:         For each  $(\mathbf{x}_i, y_i) \in \mathcal{B}$ , sample a noise  $\mathbf{z}_i$  from  $U(-1, 1)$ .
  - 5:         Use the sampled noises to approximately calculate the gradient of Eq. (A.3) w.r.t.  $\theta, \phi$ , and  $\psi$ , and obtain  $\mathbf{g}_\theta, \mathbf{g}_\phi$ , and  $\mathbf{g}_\psi$ .
  - 6:         Update  $\theta$  by:  $\theta \leftarrow \theta - \eta_\theta \cdot \mathbf{g}_\theta$ .
  - 7:         Update  $\phi$  by:  $\phi \leftarrow \phi + \eta_\phi \cdot \mathbf{g}_\phi$ .
  - 8:         Update  $\psi$  by:  $\psi \leftarrow \psi + \eta_\psi \cdot \mathbf{g}_\psi$ .
  - 9:     **end for**
  - 10: **end for**
- 

*Proof.* Due to the diagonal covariance matrix, each dimension of  $p_{\phi_i}(\delta_i)$  is independent. Thus we only consider one dimension of  $\delta_i$ . For clarity, we denote  $\mu_i^{(j)}, \sigma_i^{(j)}, \mathbf{r}^{(j)}, \mathbf{u}_i^{(j)}$ , and  $\delta_i^{(j)}$  as  $\mu, \sigma, r, u$ , and  $\delta$ , respectively. The probability density function of  $\delta$  is (see Appendix B.2 for details)

$$\begin{aligned} p(\delta) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\left(\frac{1}{2} \log \frac{\epsilon+\delta}{\epsilon-\delta} - \mu\right)^2}{2\sigma^2}\right) \cdot \frac{\epsilon}{\epsilon^2 - \delta^2} \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{r^2}{2}\right) \cdot \frac{1}{1 - \tanh(\mu + \sigma r)^2} \cdot \frac{1}{\epsilon}. \end{aligned}$$

By calculation, we have

$$\begin{aligned} p(\delta) &= \frac{1}{4\sqrt{2\pi}\sigma\epsilon} \left[ \exp\left(-\frac{r^2}{2} + 2\sigma r + 2\mu\right) + 2 \exp\left(-\frac{r^2}{2}\right) + \exp\left(-\frac{r^2}{2} - 2\sigma r - 2\mu\right) \right] \\ &\leq \frac{1}{4\sqrt{2\pi}\sigma\epsilon} \left[ \exp(2\sigma^2 + 2\mu) + 2 + \exp(2\sigma^2 - 2\mu) \right] \\ &\leq \frac{1}{4\sqrt{2\pi}\kappa_\sigma^{lo}\epsilon} \left[ 2 \exp(2(\kappa_\sigma^{up})^2) + 2\kappa_\mu + 2 \right]. \end{aligned}$$

Hence,  $p(\delta)$  is bounded. And the probability density function  $p_{\phi_i}(\delta_i)$  is also bounded since it equals to the product of  $p(\delta)$  across all dimensions.

We next prove  $p(\delta)$  is Lipschitz continuous at  $\delta \in (-\epsilon, \epsilon)$ . By calculating the derivative of  $p(\delta)$ , we have

$$\begin{aligned} p'(\delta) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\left(\frac{1}{2} \log \frac{\epsilon+\delta}{\epsilon-\delta} - \mu\right)^2}{2\sigma^2}\right) \cdot \left[ \frac{2\epsilon\delta}{(\epsilon^2 - \delta^2)^2} + \frac{\frac{1}{2} \log \frac{\epsilon+\delta}{\epsilon-\delta} - \mu}{\sigma^2} \cdot \left(\frac{\epsilon}{\epsilon^2 - \delta^2}\right)^2 \right] \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{r^2}{2}\right) \cdot \left[ \frac{2 \tanh(\mu + \sigma r)}{\epsilon^2(1 - \tanh(\mu + \sigma r)^2)^2} + \frac{r}{\sigma\epsilon^2(1 - \tanh(\mu + \sigma r)^2)^2} \right]. \end{aligned}$$

Note that although  $p'(\delta)$  has a more complicated form, the quadratic term inside  $\exp$  is still  $-\frac{r^2}{2}$ . Hence,  $p'(\delta)$  can also be bounded by a constant. Then  $p(\delta)$  as well as  $p_{\phi_i}(\delta_i)$  are Lipschitz continuous. The Lipschitz constant only concerns with  $\epsilon, \kappa_\mu, \kappa_\sigma^{lo}$ , and  $\kappa_\sigma^{up}$ . Hence, the set of explicit distributions in  $\mathcal{P}$  with a common Lipschitz constant is equicontinuous.

Combining the results, we prove that the probability density functions of the set of explicit adversarial distributions defined in Eq. (8) are bound and equicontinuous, which satisfies Assumption 2.  $\square$

**Remark 2.** Assumption 2 is used to make the search space  $\mathcal{P}$  of the inner problem in ADT compact, as can be seen in Lemma 1. However, it is a sufficient but not necessary condition of making  $\mathcal{P}$  compact. For example, if  $\mathcal{P}$  only contains Delta distributions, ADT degenerates to the AT formulation in Eq. (1) and  $\mathcal{P}$  can be represented by  $\mathcal{S}$ . In this case, it is easy to see that Assumption 2 is not satisfied but the search space of the inner problem is also compact.

**Theorem 1.** Suppose Assumptions 1 and 2 hold. We define  $\rho(\boldsymbol{\theta}) = \max_{p(\boldsymbol{\delta}_i) \in \mathcal{P}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta})$ , and  $\mathcal{P}^*(\boldsymbol{\theta}) = \{p(\boldsymbol{\delta}_i) \in \mathcal{P} : \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}) = \rho(\boldsymbol{\theta})\}$ . Then  $\rho(\boldsymbol{\theta})$  is directionally differentiable, and its directional derivative along the direction  $\mathbf{v}$  satisfies

$$\rho'(\boldsymbol{\theta}; \mathbf{v}) = \sup_{p(\boldsymbol{\delta}_i) \in \mathcal{P}^*(\boldsymbol{\theta})} \mathbf{v}^\top \nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}). \quad (\text{B.1})$$

Particularly, when  $\mathcal{P}^*(\boldsymbol{\theta}) = \{p^*(\boldsymbol{\delta}_i)\}$  only contains one maximizer,  $\rho(\boldsymbol{\theta})$  is differentiable at  $\boldsymbol{\theta}$  and

$$\nabla_{\boldsymbol{\theta}} \rho(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathcal{J}(p^*(\boldsymbol{\delta}_i), \boldsymbol{\theta}). \quad (\text{B.2})$$

*Proof.* Recall that  $\mathcal{P}$  is a set of distributions, which can be expressed by their probability density functions. The support of these functions is contained in  $\mathcal{S}$  and these functions are equicontinuous by Assumption 2.  $\mathcal{S} = \{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_\infty \leq \epsilon\}$  is the allowed perturbation set. The Euclidean distance  $\ell_2$  defines a metric on  $\mathcal{S}$ . We let

$$\mathcal{C}(\mathcal{S}, \mathbb{R}) = \{h : \mathcal{S} \rightarrow \mathbb{R} | h \text{ is continuous}\}$$

be the collection of all continuous functions from  $\mathcal{S}$  to  $\mathbb{R}$ . Then  $\mathcal{P}$  is a subset of  $\mathcal{C}(\mathcal{S}, \mathbb{R})$ . We let

$$d_{\mathcal{C}}(p, q) = \max_{\boldsymbol{\delta} \in \mathcal{S}} |p(\boldsymbol{\delta}) - q(\boldsymbol{\delta})|$$

for all  $p, q \in \mathcal{C}(\mathcal{S}, \mathbb{R})$  be a metric on  $\mathcal{C}(\mathcal{S}, \mathbb{R})$ . Then we can see that  $(\mathcal{C}(\mathcal{S}, \mathbb{R}), d_{\mathcal{C}})$  is a metric space.

We state the following lemma to prove that  $\mathcal{P}$  is compact.

**Lemma 1. (Arzelà-Ascoli's Theorem)** Let  $(X, d_X)$  be a compact metric space. A subset  $\mathcal{K}$  of  $\mathcal{C}(X, \mathbb{R})$  is compact if and only if it is closed, bounded, and equicontinuous.

Since  $(\mathcal{S}, \ell_2)$  is a compact metric space, and  $\mathcal{P}$  is closed, bounded, and equicontinuous given by Assumption 2, we can see that  $\mathcal{P}$  is compact by Lemma 1.

We next need to prove that the loss function  $\mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta})$  is continuously differentiable w.r.t. both  $p(\boldsymbol{\delta}_i)$  and  $\boldsymbol{\theta}$ , i.e., the gradient  $\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta})$  is joint continuous on  $\mathcal{P} \times \mathbb{R}^m$ , where  $m$  is the dimension of  $\boldsymbol{\theta}$ .

To prove it, we first define a new metric on  $\mathcal{P} \times \mathbb{R}^m$  as

$$d_{mix}((p_1, \boldsymbol{\theta}_1), (p_2, \boldsymbol{\theta}_2)) = d_{\mathcal{C}}(p_1, p_2) + \ell_2(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2).$$

Then  $(\mathcal{P} \times \mathbb{R}^m, d_{mix})$  is a new metric space.

By definition, given a point  $(p_0, \boldsymbol{\theta}_0) \in \mathcal{P} \times \mathbb{R}^m$ , if for each  $\tau > 0$ , there is a  $\gamma > 0$ , such that

$$\ell_2(\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}), \nabla_{\boldsymbol{\theta}} \mathcal{J}(p_0(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)) < \tau$$

whenever  $d_{mix}((p, \boldsymbol{\theta}), (p_0, \boldsymbol{\theta}_0)) < \gamma$ , then the function  $\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta})$  is continuous at  $(p_0, \boldsymbol{\theta}_0)$ . If for all points in  $\mathcal{P} \times \mathbb{R}^m$ , the function is continuous, then  $\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta})$  is continuous on  $\mathcal{P} \times \mathbb{R}^m$ .

To show that, we first have

$$\begin{aligned} & \ell_2(\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}), \nabla_{\boldsymbol{\theta}} \mathcal{J}(p_0(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)) \\ & \leq \ell_2(\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}), \nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)) + \ell_2(\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0), \nabla_{\boldsymbol{\theta}} \mathcal{J}(p_0(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)). \end{aligned} \quad (\text{B.3})$$

We already have that the loss function  $\mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta})$  is continuously differentiable w.r.t.  $\boldsymbol{\theta}$  by Assumption 1. Then given  $\frac{\tau}{2}$ , there is a  $\gamma_1$ , such that

$$\ell_2(\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}), \nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)) < \frac{\tau}{2}$$

whenever  $\ell_2(\boldsymbol{\theta}, \boldsymbol{\theta}_0) < \gamma_1$ .

For the second term of the RHS of Eq. (B.3), we have

$$\begin{aligned} \ell_2(\nabla_{\boldsymbol{\theta}} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0), \nabla_{\boldsymbol{\theta}} \mathcal{J}(p_0(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)) &= \|\nabla_{\boldsymbol{\theta}} (\mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0) - \mathcal{J}(p_0(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0))\|_2 \\ &= \left\| \int_{\mathcal{S}} (p(\boldsymbol{\delta}_i) - p_0(\boldsymbol{\delta}_i)) \nabla_{\boldsymbol{\theta}} \mathcal{L}(f_{\boldsymbol{\theta}}(\mathbf{x}_i + \boldsymbol{\delta}_i), y_i) d\boldsymbol{\delta}_i \right\|_2 \\ &\leq d_{\mathcal{C}}(p, p_0) \cdot \int_{\mathcal{S}} \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(f_{\boldsymbol{\theta}}(\mathbf{x}_i + \boldsymbol{\delta}_i), y_i)\|_2 d\boldsymbol{\delta}_i. \end{aligned}$$

Therefore, for the given  $\frac{\tau}{2}$ , there is also a  $\gamma_2$  which equals to

$$\gamma_2 = \frac{\tau}{2 \int_{\mathcal{S}} \|\nabla_{\theta} \mathcal{L}(f_{\theta}(\mathbf{x}_i + \boldsymbol{\delta}_i), y_i)\|_2 d\boldsymbol{\delta}_i},$$

such that

$$\ell_2(\nabla_{\theta} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0), \nabla_{\theta} \mathcal{J}(p_0(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)) < \frac{\tau}{2}$$

whenever  $d_{\mathcal{C}}(p, p_0) < \gamma_2$ .

Combining the results, for a given  $\tau > 0$ , we can set  $\gamma = \gamma_1 + \gamma_2$ , such that

$$\ell_2(\nabla_{\theta} \mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta}), \nabla_{\theta} \mathcal{J}(p_0(\boldsymbol{\delta}_i), \boldsymbol{\theta}_0)) < \tau$$

whenever  $d_{mix}((p, \boldsymbol{\theta}), (p_0, \boldsymbol{\theta}_0)) < \gamma$ . Thus we have proven that the loss function  $\mathcal{J}(p(\boldsymbol{\delta}_i), \boldsymbol{\theta})$  is continuously differentiable w.r.t. both  $p(\boldsymbol{\delta}_i)$  and  $\boldsymbol{\theta}$ .

Given the above results, we can directly apply Danskin's theorem [4] to prove Theorem 1. We state the Danskin's theorem in the following lemma.

**Lemma 2. (Danskin's Theorem)** *Let  $\mathcal{Q}$  be a nonempty compact topological space and  $h : \mathcal{Q} \times \mathbb{R}^m \rightarrow \mathbb{R}$  be a function satisfying that  $h(q, \cdot)$  is differentiable for every  $q \in \mathcal{Q}$  and  $\nabla_{\theta} h(q, \boldsymbol{\theta})$  is continuous on  $\mathcal{Q} \times \mathbb{R}^m$ . We define  $\Psi(\boldsymbol{\theta}) = \max_{q \in \mathcal{Q}} h(q, \boldsymbol{\theta})$ , and  $\mathcal{Q}^*(\boldsymbol{\theta}) = \{q \in \mathcal{Q} : h(q, \boldsymbol{\theta}) = \Psi(\boldsymbol{\theta})\}$ . Then  $\Psi(\boldsymbol{\theta})$  is directionally differentiable, and its directional derivative along the direction  $\mathbf{v}$  satisfies*

$$\Psi'(\boldsymbol{\theta}; \mathbf{v}) = \sup_{q \in \mathcal{Q}^*(\boldsymbol{\theta})} \mathbf{v}^{\top} \nabla_{\theta} h(q, \boldsymbol{\theta}).$$

Particularly, when  $\mathcal{Q}^*(\boldsymbol{\theta}) = \{q^*\}$  only contains one maximizer,  $\Psi(\boldsymbol{\theta})$  is differentiable at  $\boldsymbol{\theta}$  and

$$\nabla_{\theta} \Psi(\boldsymbol{\theta}) = \nabla_{\theta} h(q^*, \boldsymbol{\theta}).$$

If we let  $\mathcal{Q} = \mathcal{P}$  and  $h = \mathcal{J}$  in Lemma 2, we can directly prove Theorem 1. □

## B.2 Proof of Eq. (11)

The variable  $\boldsymbol{\delta}_i$  has the following sampling process

$$\boldsymbol{\delta}_i = \epsilon \cdot \tanh(\mathbf{u}_i), \quad \mathbf{u}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i^2)),$$

whose negative log density is

$$\sum_{j=1}^d \left( \frac{1}{2} (\mathbf{r}^{(j)})^2 + \frac{\log 2\pi}{2} + \log \boldsymbol{\sigma}_i^{(j)} + \log(1 - \tanh(\boldsymbol{\mu}_i^{(j)} + \boldsymbol{\sigma}_i^{(j)} \mathbf{r}^{(j)})^2) + \log \epsilon \right),$$

where the superscript  $j$  denotes the  $j$ -th element of a vector.

*Proof.* Due to the usage of the diagonal covariance matrix, each dimension in the sampled perturbation  $\boldsymbol{\delta}_i$  is independent. Thus we can simply calculate the negative log density in each dimension of  $\boldsymbol{\delta}_i$ . For clarity, we also denote  $\boldsymbol{\mu}_i^{(j)}$ ,  $\boldsymbol{\sigma}_i^{(j)}$ ,  $\mathbf{r}^{(j)}$ ,  $\mathbf{u}_i^{(j)}$ , and  $\boldsymbol{\delta}_i^{(j)}$  as  $\mu$ ,  $\sigma$ ,  $r$ ,  $u$ , and  $\delta$ , respectively. Based on the sampling procedure in Eq. (8), we have  $\delta = \epsilon \cdot \tanh(u)$  and  $u = \mu + \sigma r$ .

Note that  $r$  has density:  $p(r) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{r^2}{2})$ . Apply the *transformation of variable* approach, we have the density of  $u$  as

$$p(u) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{r^2}{2}) \cdot \left| \frac{d}{du} \left( \frac{u - \mu}{\sigma} \right) \right| = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{r^2}{2}).$$

Let  $\beta = \tanh(u)$ , then the inverse transformation is  $u = \tanh^{-1}(\beta) = \frac{1}{2} \log\left(\frac{1+\beta}{1-\beta}\right)$ , whose derivative w.r.t.  $\beta$  is  $\frac{1}{1-\beta^2}$ .

Then, by applying the *transformation of variable* approach again, we have the density of  $\beta$  as

$$p(\beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{r^2}{2}) \cdot \frac{1}{1-\beta^2} = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{r^2}{2}) \cdot \frac{1}{1-\tanh(\mu + \sigma r)^2}.$$

Therefore, the density of  $\delta$  which equals to  $\epsilon \cdot \beta$  can be derived similarly, and eventually we obtain

$$p(\delta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{r^2}{2}\right) \cdot \frac{1}{1 - \tanh(\mu + \sigma r)^2} \cdot \frac{1}{\epsilon}.$$

Consequently, the negative log density of  $p(\delta)$  is

$$-\log p(\delta) = \frac{r^2}{2} + \frac{\log 2\pi}{2} + \log \sigma + \log(1 - \tanh(\mu + \sigma r)^2) + \log \epsilon.$$

Sum over all of the dimensions and we complete the proof of Eq. (11).  $\square$

### B.3 Proof of Eq. (A.2)

Given an example  $\mathbf{x}_i$ , we define an implicit adversarial distribution  $p_\phi(\delta_i|\mathbf{x}_i)$  in the form of  $\delta_i = g_\phi(\mathbf{z}; \mathbf{x}_i)$ ,  $\mathbf{z} \sim p(\mathbf{z})$ , where  $g_\phi$  denotes a  $\phi$ -parameterized generator network. Then we can maximize the following variational lower bound to maximize the entropy of  $p_\phi(\delta_i|\mathbf{x}_i)$ , as

$$\mathcal{H}(p_\phi(\delta_i|\mathbf{x}_i)) \geq \mathcal{U}(q) = \mathbb{E}_{p(\mathbf{z})} \log q(\mathbf{z}|g_\phi(\mathbf{z}; \mathbf{x}_i)) + c$$

where  $c$  is a constant and  $q(\cdot)$  is an introduced variational distribution.

*Proof.* We mainly follow [2] to provide the proof. Typically, we can view the Dirac generation distribution  $p_\phi(\delta_i|\mathbf{x}_i, \mathbf{z})$  as a peaked Gaussian with a fixed, diagonal covariance, then it will have a constant entropy. Considering  $\mathbf{x}_i$  as a given condition, we can simply rewrite the generation distribution as  $p_{\phi,i}(\delta_i|\mathbf{z})$ . Then we can define the joint distribution over  $\delta_i$  and  $\mathbf{z}$  as  $p_{\phi,i}(\delta_i, \mathbf{z}) = p_{\phi,i}(\delta_i|\mathbf{z})p_{\phi,i}(\mathbf{z})$ .  $p_{\phi,i}(\mathbf{z}) = p(\mathbf{z})$  is simply a predefined prior with a constant entropy. Then, we can further define the marginal  $p_{\phi,i}(\delta_i)$  whose entropy is of our interest and the posterior  $p_{\phi,i}(\mathbf{z}|\delta_i)$ . Consider the mutual information between  $\delta_i$  and  $\mathbf{z}$

$$\mathcal{I}(p_{\phi,i}(\delta_i); p_{\phi,i}(\mathbf{z})) = \mathcal{H}(p_{\phi,i}(\delta_i)) - \mathcal{H}(p_{\phi,i}(\delta_i|\mathbf{z})) = \mathcal{H}(p_{\phi,i}(\mathbf{z})) - \mathcal{H}(p_{\phi,i}(\mathbf{z}|\delta_i)).$$

Thus, we can calculate the entropy of  $\delta_i$  as

$$\mathcal{H}(p_{\phi,i}(\delta_i)) = \mathcal{H}(p_{\phi,i}(\mathbf{z})) - \mathcal{H}(p_{\phi,i}(\mathbf{z}|\delta_i)) + \mathcal{H}(p_{\phi,i}(\delta_i|\mathbf{z})).$$

As stated, the first term and the last term are constant w.r.t. the parameter  $\phi$ . Therefore, maximizing  $\mathcal{H}(p_{\phi,i}(\delta_i))$  corresponds to maximizing the negative conditional entropy

$$-\mathcal{H}(p_{\phi,i}(\mathbf{z}|\delta_i)) = \mathbb{E}_{\delta_i \sim p_{\phi,i}(\delta_i)} [\mathbb{E}_{\mathbf{z} \sim p_{\phi,i}(\mathbf{z}|\delta_i)} [\log p_{\phi,i}(\mathbf{z}|\delta_i)]].$$

We still cannot optimize this as we have no access to the posterior. As an alternative, we resort to the variational inference technique to tackle this problem. We introduce a variational distribution  $q(\mathbf{z}|\delta_i)$  to approximate the true posterior, and derive the following lower bound

$$\begin{aligned} -\mathcal{H}(p_{\phi,i}(\mathbf{z}|\delta_i)) &= \mathbb{E}_{\delta_i \sim p_{\phi,i}(\delta_i)} [\mathbb{E}_{\mathbf{z} \sim p_{\phi,i}(\mathbf{z}|\delta_i)} [\log q(\mathbf{z}|\delta_i)]] + \mathcal{D}_{KL}(p_{\phi,i}(\mathbf{z}|\delta_i) || q(\mathbf{z}|\delta_i)) \\ &\geq \mathbb{E}_{\delta_i \sim p_{\phi,i}(\delta_i)} [\mathbb{E}_{\mathbf{z} \sim p_{\phi,i}(\mathbf{z}|\delta_i)} [\log q(\mathbf{z}|\delta_i)]] \\ &= \mathbb{E}_{\mathbf{z}, \delta_i \sim p_{\phi,i}(\mathbf{z}, \delta_i)} [\log q(\mathbf{z}|\delta_i)] \\ &= \underbrace{\mathbb{E}_{\mathbf{z} \sim p_{\phi,i}(\mathbf{z})} [\mathbb{E}_{\delta_i \sim p_{\phi,i}(\delta_i|\mathbf{z})} [\log q(\mathbf{z}|\delta_i)]]}_{\mathcal{U}'(q)}, \end{aligned}$$

where  $\mathcal{D}_{KL}$  represents the Kullback–Leibler divergence between distributions. Note that  $p_{\phi,i}(\mathbf{z}) = p(\mathbf{z})$  is a prior and  $p_{\phi,i}(\delta_i|\mathbf{z}) = p_\phi(\delta_i|\mathbf{x}_i, \mathbf{z})$  is Dirac distribution located at  $\delta_i = g_\phi(\mathbf{z}; \mathbf{x}_i)$ . Thus, we can write the lower bound of the entropy  $\mathcal{U}(q)$  as

$$\mathcal{H}(p_\phi(\delta_i|\mathbf{x}_i)) \geq \mathcal{U}(q) = \mathcal{U}'(q) + c = \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \log q(\mathbf{z}|g_\phi(\mathbf{z}; \mathbf{x}_i)) + c,$$

which can be optimized effectively via Monte Carlo integration and standard back-propagation. Then we finish the proof of Eq. (A.2).  $\square$

## C Detailed experimental settings

We provide the detailed experimental settings in this section. All of the experiments are conducted on NVIDIA 2080 Ti GPUs. The source code of ADT is available at <https://github.com/dongyp13/Adversarial-Distributional-Training>.

Table 7: The network architectures used for the generators.

In ADT <sub>EXP-AM</sub>	In ADT <sub>IMP-AM</sub>
input	$\mathbf{z}$
$256 \times 3 \times 3$ conv	256-dim fc layer
Residual block, 512 filters	1024-dim fc layer
Residual block, 512 filters	reshape to $1 \times 32 \times 32$
Residual block, 512 filters	concat with input
$6 \times 3 \times 3$ conv	$256 \times 3 \times 3$ conv
	Residual block, 512 filters
	Residual block, 512 filters
	Residual block, 512 filters
	$3 \times 3 \times 3$ conv

Table 8: The network architecture used for instantiating the variational distribution  $q$  in ADT<sub>IMP-AM</sub>.

Layers
input
$32 \times 5 \times 5$ , stride 1
$64 \times 4 \times 4$ , stride 2
$128 \times 4 \times 4$ , stride 1
$256 \times 4 \times 4$ , stride 2
Global average pooling
$128 \times 1 \times 1$ , stride 1

## C.1 Datasets

We choose the CIFAR-10 [10], CIFAR-100 [10], and SVHN [14] datasets to conduct the experiments. CIFAR consists of a training set of 50,000 and a test set of 10,000 color images of resolution  $32 \times 32$  with 10 classes in CIFAR-10 and 100 classes in CIFAR-100. SVHN is a 10-class house number classification dataset with 73,257 training images and 26,032 test images. During training, we perform standard data augmentation (i.e., horizontal flips and random crops from images with 4 pixels padded on each side) on CIFAR-10 and CIFAR-100, and use no data augmentation on SVHN. We do not use any data augmentation during testing.

## C.2 Network architectures

For the generator network in ADT<sub>EXP-AM</sub> and ADT<sub>IMP-AM</sub>, we adopt a popular image-to-image architecture which has shown promise in neural style transfer and super-resolution [7, 18]. The network contains 3 residual blocks [5], with two extra convolutions at the beginning and the end. All convolutions in the generator have stride 1, and are immediately followed by batch normalization [6] and ReLU activation.

As found by [1], taking only the natural images as inputs to the generator network can lead to poor results. And they suggest to input the classifier’s gradients as well. Based on this finding, we calculate the gradient of the loss function at the natural input  $\mathbf{g}_i^1 = \nabla_{\mathbf{x}} \mathcal{L}(f_{\theta}(\mathbf{x}_i), y_i)$ , as well as the gradient of the loss function at the FGSM adversarial example  $\mathbf{g}_i^2 = \nabla_{\mathbf{x}} \mathcal{L}(f_{\theta}(\mathbf{x}_i + \delta_i^{\text{FGSM}}), y_i)$ , where  $\delta_i^{\text{FGSM}} = \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}(f_{\theta}(\mathbf{x}_i), y_i))$ , and then input  $[\mathbf{x}_i, \mathbf{g}_i^1, \mathbf{g}_i^2]$  to the generator network.

In ADT<sub>EXP-AM</sub>, the generator has 6 output channels to deliver the parameters (i.e., mean and standard derivation) of the explicit adversarial distributions. In ADT<sub>IMP-AM</sub>, for each input we sample a 64-dim i.i.d.  $\mathbf{z}$  from a uniform distribution  $U(-1, 1)$ , which is encoded with 2 fully connected (FC) layers and then fed into the generator along with the input image and gradients.

We elaborate the architectures of the generator networks in Table 7, and the architecture of  $q$  in ADT<sub>IMP-AM</sub> in Table 8. In these tables, “ $C \times H \times W$ ” means a convolutional layer with  $C$  filters size  $H \times W$ , which is followed by batch normalization [6] and a ReLU nonlinearity (or LeakyReLU for layers in Table 8), except the last layers in the architectures. We use the residual block design in [5], which is composed of two  $3 \times 3$  convolutions and a residual connection.

Table 9: Classification accuracy of the three proposed methods and baselines on CIFAR-100 and SVHN under white-box attacks. We mark the best results for each attack and the overall results that outperform the baselines in **bold**, and the overall best result in **blue**.

Model	$\mathcal{A}_{\text{nat}}$	FGSM	PGD-20	PGD-100	MIM	C&W	FeaAttack	$\mathcal{A}_{\text{rob}}$
CIFAR-100, $\epsilon = 8/255$								
Standard	<b>78.59%</b>	8.73%	0.02%	0.01%	0.02%	0.00%	0.00%	0.00%
AT <sub>PGD</sub>	61.45%	30.78%	25.71%	25.40%	26.60%	25.80%	33.95%	24.49%
ADT <sub>EXP</sub>	62.70%	34.22%	28.96%	<b>28.60%</b>	<b>29.83%</b>	<b>28.99%</b>	<b>35.07%</b>	<b>27.13%</b>
ADT <sub>EXP-AM</sub>	62.84%	36.28%	29.01%	28.46%	29.68%	28.78%	34.91%	<b>26.87%</b>
ADT <sub>IMP-AM</sub>	64.07%	<b>39.39%</b>	<b>29.40%</b>	28.43%	29.64%	28.76%	35.00%	<b>26.80%</b>
SVHN, $\epsilon = 4/255$								
Standard	<b>96.12%</b>	39.05%	3.64%	2.95%	4.08%	3.91%	2.14%	2.14%
AT <sub>PGD</sub>	95.07%	82.19%	74.22%	73.79%	74.56%	74.77%	73.51%	73.38%
ADT <sub>EXP</sub>	95.70%	86.72%	<b>77.01%</b>	<b>76.62%</b>	<b>77.18%</b>	<b>77.50%</b>	<b>75.64%</b>	<b>75.55%</b>
ADT <sub>EXP-AM</sub>	95.67%	85.24%	76.12%	75.58%	76.63%	76.70%	75.20%	<b>75.00%</b>
ADT <sub>IMP-AM</sub>	95.62%	<b>86.73%</b>	75.61%	74.85%	75.91%	76.12%	74.24%	<b>74.13%</b>

### C.3 Training details

The classifier is trained using SGD with momentum 0.9, weight decay  $2 \times 10^{-4}$ , and batch size 64. The initial learning rate is 0.1, which is reduced to 0.01 in the 75-th epoch. We stop training after 76 epochs. For ADT<sub>EXP</sub>, we adopt Adam [9] for optimizing the distribution parameters  $\phi_i$ . We set the learning rate for  $\phi_i$  as 0.3, the momentum as (0.0, 0.0), the number of optimization steps as  $T = 7$ , and the number of MC samples to estimate the gradient in each step as  $k = 5$ . For ADT<sub>EXP-AM</sub> and ADT<sub>IMP-AM</sub>, we use only one MC sample for gradient estimation and use Adam with momentum (0.5, 0.999) and learning rate  $2 \times 10^{-4}$  to optimize the parameter  $\phi$  of the generator network. We also adopt Adam with learning rate  $2 \times 10^{-4}$  to optimize the parameter  $\psi$  of the introduced variational in ADT<sub>IMP-AM</sub>.

### C.4 Baselines

Our primary baselines include: 1) standard training on the clean images (*Standard*); 2) adversarial training on the PGD adversarial examples (AT<sub>PGD</sub>) [13]. Standard and AT<sub>PGD</sub> are trained with the same configurations specified above. For training AT<sub>PGD</sub>, we perform PGD with  $T = 7$  steps, and step size  $\alpha = \epsilon/4$ , which are the same as in [13]. On CIFAR-10, we incorporate several additional baselines, including: 1) the pretrained AT<sub>PGD</sub> model (AT<sub>PGD</sub><sup>†</sup>) released by [13]; 2) adversarial training on the targeted FGSM adversarial examples (AT<sub>FGSM</sub>) [11]; 3) adversarial logit pairing (ALP) [8]; and 4) feature scattering-based adversarial training (FeaScatter) [16]. We implement AT<sub>FGSM</sub> and ALP by ourselves using the same training configuration specified above and use the pretrained model of FeaScatter. Note that all of these models have the same network architecture for a fair comparison.

### C.5 A feature attack for white-box evaluation

We incorporate a feature attack (FeaAttack) [12] for white-box robustness evaluation in this paper. The algorithm of FeaAttack is introduced below. Given a natural input  $\mathbf{x}$ , FeaAttack first finds a target image  $\mathbf{x}'$  belonging to a different class. It minimizes the cosine similarity between the feature representations of the adversarial example and  $\mathbf{x}'$  as

$$\delta^* = \arg \min_{\delta \in \mathcal{S}} \mathcal{L}_{\text{cos}}(f'_{\theta}(\mathbf{x} + \delta), f'_{\theta}(\mathbf{x}')),$$

where  $f'_{\theta}(\cdot)$  returns the feature representation before the global average pooling layer for an input, and  $\mathcal{L}_{\text{cos}}$  is the cosine similarity between two features. FeaAttack solves this objective function by

$$\delta^{t+1} = \Pi_{\mathcal{S}}(\delta^t - \alpha \cdot \text{sign}(\nabla_{\mathbf{x}} \mathcal{L}_{\text{cos}}(f'_{\theta}(\mathbf{x} + \delta^t), f'_{\theta}(\mathbf{x}')))).$$

$\delta^0$  is initialized uniformly in  $\mathcal{S}$ . In our experiments, we set  $\alpha = \epsilon/8$  and the number of optimization steps as 50. For each natural input, we randomly select 200 target images to conduct 200 attacks, and report a successful attack when one of them can cause misclassification of the model.



Table 10: Classification accuracy of TRADES and the three ADT-based methods trained with the TRADES loss on CIFAR-10 under white-box attacks with  $\epsilon = 8/255$ . We mark the best results for each attack and the overall results that outperform the baselines in **bold**, and the overall best result in **blue**.

Model	$\beta$	$\mathcal{A}_{\text{nat}}$	FGSM	PGD-20	PGD-100	MIM	C&W	FeaAttack	$\mathcal{A}_{\text{rob}}$
TRADES	1.0	87.99%	57.67%	51.08%	48.41%	53.32%	49.29%	51.07%	47.75%
ADT <sub>EXP</sub>	1.0	<b>89.74%</b>	59.47%	52.39%	49.88%	54.74%	50.75%	51.29%	<b>49.05%</b>
ADT <sub>EXP-AM</sub>	1.0	88.86%	62.89%	<b>54.44%</b>	<b>51.66%</b>	<b>56.09%</b>	<b>52.33%</b>	<b>54.61%</b>	<b>50.78%</b>
ADT <sub>IMP-AM</sub>	1.0	88.80%	<b>68.35%</b>	54.22%	51.09%	54.95%	51.84%	54.19%	<b>50.14%</b>
TRADES	6.0	84.02%	60.08%	56.06%	54.49%	57.27%	53.62%	55.18%	52.64%
ADT <sub>EXP</sub>	6.0	84.66%	61.72%	57.71%	<b>56.17%</b>	<b>58.74%</b>	<b>55.16%</b>	56.65%	<b>54.21%</b>
ADT <sub>EXP-AM</sub>	6.0	84.85%	66.09%	57.67%	55.73%	58.38%	54.79%	58.94%	<b>54.09%</b>
ADT <sub>IMP-AM</sub>	6.0	<b>84.96%</b>	<b>68.34%</b>	<b>57.82%</b>	55.45%	58.58%	54.36%	<b>59.01%</b>	<b>53.66%</b>

## D Supplementary experimental results

We provide more experimental results in this section.

### D.1 Full results on CIFAR-100 and SVHN

We provide the full experimental results of Standard, AT<sub>PGD</sub>, ADT<sub>EXP</sub>, ADT<sub>EXP-AM</sub>, and ADT<sub>IMP-AM</sub> under all adopted white-box attacks on CIFAR-100 and SVHN in Table 9.

### D.2 Full results on TRADES

In TRADES [17], the minimax optimization problem is formulated as

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \left\{ \mathcal{L}(f_{\theta}(\mathbf{x}_i), y_i) + \beta \cdot \max_{\delta_i \in \mathcal{S}} \mathcal{D}_{\text{KL}}(f_{\theta}(\mathbf{x}_i + \delta_i), f_{\theta}(\mathbf{x}_i)) \right\},$$

where  $\beta$  is a hyperparameter balancing the trade-off between natural and robust accuracy. The full experimental results of TRADES and the three variants of ADT when integrated with the TRADES loss are shown in Table 10. We evaluate their performance by all adopted white-box attacks and report the worst-case robustness as in Eq. (12).

### D.3 Convergence of learning the explicit adversarial distributions

We study the convergence of the explicit adversarial distributions introduced in Sec. 3.1 by attacking AT<sub>PGD</sub> and ADT<sub>EXP</sub> with varying iterations. We set the learning rate of  $\phi_i$  as 0.3, the momentum as (0.0, 0.0), the number of MC samples to estimate the gradient in each step as  $k = 10$ , and vary the attack iterations from 0 to 100. We show the classification loss and accuracy in Fig. 6. Learning the explicit adversarial distributions can converge soon within a few iterations.

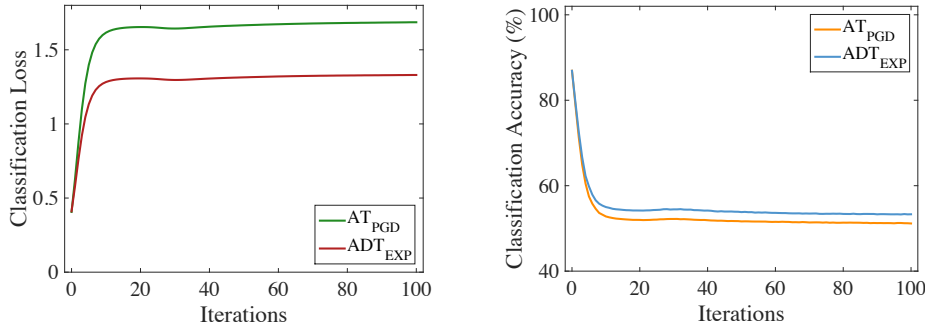


Figure 6: Classification loss (i.e., cross-entropy loss) and accuracy (%) of AT<sub>PGD</sub> and ADT<sub>EXP</sub> under the explicit adversarial distributions attack with different attack iterations.

#### D.4 Training time

We provide the one-epoch training time of Standard,  $AT_{PGD}$ ,  $ADT_{EXP}$ ,  $ADT_{EXP-AM}$ , and  $ADT_{IMP-AM}$  on CIFAR-10 in Fig. 7. As can be seen,  $ADT_{EXP}$  is nearly  $5\times$  slower than  $AT_{PGD}$  since we use  $k = 5$  MC samples to estimate the gradient w.r.t. the distribution parameters in each step. Nevertheless, by amortizing the adversarial distributions,  $ADT_{EXP-AM}$  and  $ADT_{IMP-AM}$  are much faster than  $ADT_{EXP}$ , and nearly  $2\times$  faster than  $AT_{PGD}$ .

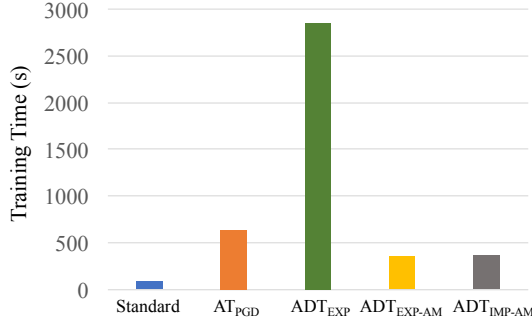


Figure 7: The training time (s) for one epoch of Standard,  $AT_{PGD}$ ,  $ADT_{EXP}$ ,  $ADT_{EXP-AM}$ , and  $ADT_{IMP-AM}$  on CIFAR-10.

#### D.5 Comparison with Chen et al. [1]

We further compare ADT with the L2L framework in [1]. Their method is similar to ours in the sense that they also adopt a generator network to produce adversarial examples, and perform adversarial training on those generated adversarial examples. The essential different between our methods and theirs is that we propose an adversarial distributional training framework to learn the distributions of adversarial perturbations, while their method is a variant of the vanilla adversarial training with a different approach to solving the inner maximization.

Since the source code is not provided by Chen et al. [1], we tried to reproduce their reported results with the same training configuration specified in their paper, but we failed. Therefore, we adopt the same configuration as in ADT for training the L2L model. Table 11 shows the results of L2L,  $ADT_{EXP-AM}$ , and  $ADT_{IMP-AM}$ , which use the same classifier architecture and generator network. Our ADT-based methods outperform L2L in most cases, showing the advantages of learning the distributions of adversarial perturbations upon finding a single adversarial example.

Table 11: Classification Accuracy of L2L [1],  $ADT_{EXP-AM}$ , and  $ADT_{IMP-AM}$  on CIFAR-10 under white-box attacks with  $\epsilon = 8/255$ .

Model	L2L	$ADT_{EXP-AM}$	$ADT_{IMP-AM}$
$\mathcal{A}_{nat}$	<b>88.15%</b>	87.82%	88.00%
FGSM	<b>65.50%</b>	62.42%	64.89%
PGD-20	48.55%	51.95%	<b>52.28%</b>
PGD-100	47.14%	<b>51.26%</b>	51.23%
MIM	49.03%	<b>52.99%</b>	52.64%
C&W	49.22%	51.75%	<b>52.65%</b>

## References

- [1] Zhehui Chen, Haoming Jiang, Yuyang Shi, Bo Dai, and Tuo Zhao. Learning to defense by learning to attack. *arXiv preprint arXiv:1811.01213*, 2018.
- [2] Zihang Dai, Amjad Almahairi, Philip Bachman, Eduard Hovy, and Aaron Courville. Calibrating energy-based generative adversarial networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [3] Zihang Dai, Zhilin Yang, Fan Yang, William W Cohen, and Russ R Salakhutdinov. Good semi-supervised learning that requires a bad gan. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

- [4] John M Danskin. *The theory of max-min and its application to weapons allocation problems*, volume 5. Springer Science & Business Media, 2012.
- [5] 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 (CVPR)*, 2016.
- [6] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, 2015.
- [7] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision (ECCV)*, 2016.
- [8] Harini Kannan, Alexey Kurakin, and Ian Goodfellow. Adversarial logit pairing. *arXiv preprint arXiv:1803.06373*, 2018.
- [9] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- [10] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [11] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *International Conference on Learning Representations (ICLR)*, 2017.
- [12] Daquan Lin. <https://github.com/Line290/FeatureAttack>.
- [13] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [14] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [15] Jiaxin Shi, Shengyang Sun, and Jun Zhu. A spectral approach to gradient estimation for implicit distributions. In *International Conference on Machine Learning (ICML)*, 2018.
- [16] Haichao Zhang and Jianyu Wang. Defense against adversarial attacks using feature scattering-based adversarial training. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [17] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning (ICML)*, 2019.
- [18] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.