

---

# Can Adversarial Training Be Manipulated By Non-Robust Features?

---

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

1 Adversarial training, originally designed to resist test-time adversarial examples,  
2 has shown to be promising in mitigating *training-time availability attacks*. This  
3 defense ability, however, is challenged in this paper. We identify a novel threat  
4 model named *stability attacks*, which aims to hinder *robust* availability by slightly  
5 manipulating the training data. Under this threat, we show that adversarial training  
6 using a conventional defense budget  $\epsilon$  provably fails to provide test robustness  
7 in a simple statistical setting, where the non-robust features of the training data  
8 can be reinforced by  $\epsilon$ -bounded perturbation. Further, we analyze the necessity of  
9 enlarging the defense budget to counter stability attacks. Finally, comprehensive  
10 experiments demonstrate that stability attacks are harmful on benchmark datasets,  
11 and thus the adaptive defense is necessary to maintain robustness.

## 1 Introduction

13 Robustness to input perturbations is crucial to machine learning deployment in various applications,  
14 such as spam filtering [13] and autonomous driving [6]. One of the most popular methods for  
15 improving test robustness is *adversarial training* [39, 1]. By augmenting the training data with  
16  $\epsilon$ -bounded and on-the-fly crafted adversarial examples, adversarial training helps the learned model  
17 resist test-time perturbations [39].

18 On the other hand, machine learning systems  
19 are vulnerable to *training-time availability at-*  
20 *tacks* [3]. In particular, small perturbations applied into the training data before training suffice to degrade the overall test performance of naturally trained models [16, 28]. Fortunately, recent work has shown that adversarial training [39] is capable of mitigating this type of threat [60, 18]. In other words, even if the training data is manipulated by an adversary to maximize the test error, adversarially trained models can still achieve considerable accuracy on clean test data. However, previous work hardly inspects the test robustness of the models, which is what adversarial training was originally proposed for [23, 39]. This naturally raises the following question:

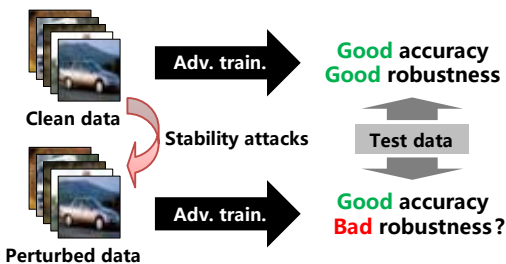


Figure 1: An illustration of stability attacks, where the training data is slightly perturbed to hinder adversarial training.

33 *Are the models adversarially trained on the manipulated data robust to test-time perturbations?*

34 In this work, we show that conventional adversarial training may fail to provide test robustness when  
 35 the training data is manipulated by an adversary, and thus an adaptive defense is necessary to resolve  
 36 this issue. Our contributions are summarized as follows:

- 37 1. We introduce a novel threat model called *stability attacks*, where an adversary aims to  
 38 degrade the overall test robustness of adversarially trained models by slightly perturbing the  
 39 training data. Figure 1 illustrates the threat of stability attacks.
- 40 2. We show that adversarial training using a conventional defense budget *provably* fails under  
 41 stability attacks in a simple statistical setting. Specifically, a defense budget of  $\epsilon$  will produce  
 42 models that are *not* robust to  $\epsilon$ -bounded adversarial examples when the training data is  
 43 hypocritically perturbed.
- 44 3. We unveil that the aforementioned vulnerability stems from the existence of the non-robust  
 45 (predictive, yet brittle) features [29] in the original training data. When the non-robust  
 46 features are reinforced by hypocritical perturbations, the conventional defense budget will  
 47 be insufficient to offset the negative impact.
- 48 4. We further show that a defense budget of  $2\epsilon$  is capable of resisting any stability attack for  
 49 adversarial training, while the budget can be reduced to  $\epsilon + \eta$  in a simple statistical setting,  
 50 where  $\eta$  is the magnitude of the non-robust features.
- 51 5. We demonstrate that stability attacks are harmful to conventional adversarial training on  
 52 benchmark datasets. In addition, our empirical study suggests that enlarging the defense  
 53 budget is essential for mitigating hypocritical perturbations.

54 To the best of our knowledge, this is the first work that studies the robustness of adversarial training  
 55 against stability attacks. Both theoretical and empirical evidences show that the conventional defense  
 56 budget  $\epsilon$  is insufficient under the threat of  $\epsilon$ -bounded training-time perturbations. Our findings suggest  
 57 that practitioners should consider a larger defense budget of no more than  $2\epsilon$  (practically, about  
 58  $1.5\epsilon \sim 1.75\epsilon$ ) to achieve a better  $\epsilon$ -robustness.

## 59 2 Threat Models

60 In this section, we formally introduce the threat model of stability attacks. We begin by revisiting the  
 61 concepts of natural risk, adversarial risk, and delusive attacks. These concepts naturally give rise to  
 62 our formulation of stability attacks.

### 63 2.1 Preliminaries

64 **Setup.** We consider a classification task with input-label pairs  $(\mathbf{x}, y)$  from an underlying distribution  
 65  $\mathcal{D}$  over  $\mathcal{X} \times [k]$ . The goal is to learn a (robust) classifier  $f : \mathcal{X} \rightarrow [k]$  that predicts a label  $y$  for a  
 66 given input  $\mathbf{x}$ .

67 **Natural training (NT).** Most learning algorithms aim to maximize the generalization performance  
 68 on unperturbed examples, i.e., natural accuracy. The goal is to minimize the natural risk defined as:

$$\mathcal{R}_{\text{nat}}(f) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathcal{L}(f(\mathbf{x}), y)]. \quad (1)$$

69 **Adversarial training (AT).** Since the risk of adversarial examples (a.k.a. evasion attacks) was  
 70 found to be unexpectedly high [5, 56], it has become increasingly important to defend the learner  
 71 against the worst-case perturbations [23, 39]. In this context, the goal is to train a model that has low  
 72 *adversarial risk* given a defense budget  $\epsilon$ :

$$\mathcal{R}_{\text{adv}}(f) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ \max_{\delta \in \Delta} \mathcal{L}(f(\mathbf{x} + \delta), y) \right], \quad (2)$$

73 where we choose  $\Delta$  to be the set of  $\epsilon$ -bounded perturbations, i.e.,  $\Delta = \{\delta \mid \|\delta\| \leq \epsilon\}$ . This choice is  
 74 the most common one in the context of adversarial examples [63]. To simplify the notation, we refer  
 75 to the robustness with respect to this set as  $\epsilon$ -robustness. It is worth noting that  $\mathcal{R}_{\text{adv}}(f) \geq \mathcal{R}_{\text{nat}}(f)$   
 76 always holds for any  $f$ , and the equation holds when  $\epsilon = 0$ .

Table 1: Comparisons between evasion attacks, delusive attacks, and stability attacks.

Threat model	Training-time perturbation	Test-time perturbation	Learning scheme	Test performance
None	✗	✗	NT	Good
Evasion attacks [5, 56, 23, 39]	✗	✓	NT AT	Bad Good
Delusive attacks [43, 16, 28, 60]	✓	✗	NT AT	Bad Good
Stability attacks (this paper)	✓	✓	AT (conventional) AT (our improved)	Bad Good

77 **Delusive attacks.** Delusive attacks, which belong to training-time availability attacks, aim to prevent  
 78 the learner from producing an accurate model by manipulating the training data “imperceptibly” [43].  
 79 Concretely, the features of the training data can be perturbed, while the labels should remain  
 80 correct [16, 53, 28, 74, 15, 60, 18, 19]. This malicious task can be formalized into the following  
 81 bi-level optimization problem:

$$\begin{aligned} & \max_{\mathcal{P} \in \mathcal{S}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathcal{L}(f_{\mathcal{P}}(\mathbf{x}), y)] \\ \text{s.t. } & f_{\mathcal{P}} \in \arg \min_f \sum_{(\mathbf{x}_i, y_i) \in \mathcal{T}} [\mathcal{L}(f(\mathbf{x}_i + \mathbf{p}_i), y_i)]. \end{aligned} \quad (3)$$

82 Here, the adversary aims to maximize the natural risk of the model  $f_{\mathcal{P}}$  (that is trained on the  
 83 manipulated training set) by applying the generated perturbations  $\mathcal{P} = \{\mathbf{p}_i\}_{i=1}^n$  into the original  
 84 training set  $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ . The commonly used feasible region is  $\mathcal{S} = \{\{\mathbf{p}_i\}_{i=1}^n \mid \|\mathbf{p}_i\| \leq \epsilon\}$ .

85 Generally, solving Equation (3) is computationally prohibitive for neural networks [60, 18]. Thus,  
 86 various heuristic methods are explored to resolve this problem. For example, a representative method  
 87 called *hypocritical perturbations* [59, 60, 28, 19] is crafted as follows:

$$\min_{\|\mathbf{p}_i\| \leq \epsilon} \mathcal{L}(f_{\text{craft}}(\mathbf{x}_i + \mathbf{p}_i), y_i), \quad (4)$$

88 where  $f_{\text{craft}}$  is called the *crafting model*, pre-trained before generating poisons. Tao et al. [60] simply  
 89 adopted a naturally trained classifier as the crafting model, while Huang et al. [28] proposed a  
 90 min-min bi-level optimization process to pre-train the crafting model. Fu et al. [19] further built their  
 91 crafting model via a min-min-max three-level optimization process, and generated their poisons by  
 92 replacing Equation (4) with a min-max bi-level objective.

93 Another representative method of delusive attacks is the *adversarial perturbation*, crafted by solving

$$\max_{\|\mathbf{p}_i\| \leq \epsilon} \mathcal{L}(f_{\text{craft}}(\mathbf{x}_i + \mathbf{p}_i), y_i). \quad (5)$$

94 Tao et al. [60] and Fowl et al. [18] independently found that applying the adversarial perturbation to  
 95 the training data is very effective at compromising naturally trained models. However, adversarial  
 96 training has proven to be promising in defending against various delusive attacks [60].

## 97 2.2 Stability Attacks

98 In contrast to delusive attacks that aim at increasing the natural risk, stability attacks attempt to  
 99 maximize the adversarial risk of the learner by slightly perturbing the training data:

$$\max_{\mathcal{P} \in \mathcal{S}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ \max_{\delta \in \Delta} \mathcal{L}(f_{\mathcal{P}}(\mathbf{x} + \delta), y) \right], \quad (6)$$

100 where  $f_{\mathcal{P}}$  denotes the victim model, which is naturally or adversarially trained on the perturbed data.  
 101 In other words, stability attacks seek to hinder the *robust* availability of the training data. Table 1  
 102 shows the comparisons among different threat models.

103 The goal of stability attacks can be immediately achieved for naturally trained models, since they  
 104 have already incurred high adversarial risk, even if the training data is clean [56]. To ease the problem

105 of high adversarial risk, adversarial training has been widely used to improve model’s adversarial  
 106 robustness [24, 12]. Hence, the main goal of stability attacks becomes to compromise the test  
 107 robustness of adversarially trained models.

108 Note that Equation (6) is a multi-level optimization problem that is not easy to solve, our next question  
 109 is how to conduct effective stability attacks against adversarial training. In the following sections, we  
 110 introduce an effective stability attack method and analyze the cost of resisting it.

111 **Justification.** This work focuses on adding bounded perturbations as small as possible. We mostly  
 112 assume that the adversary has full control of training data (instead of changing a few) by following  
 113 previous works [16, 28, 60, 17–19]. This is a realistic assumption [16, 18]. For instance, in some  
 114 applications an organization may agree to release some internal data for peer assessment, while  
 115 preventing competitors from easily building a model with high test robustness; this can be achieved  
 116 by perturbing the entire dataset via stability attacks before releasing. Moreover, this assumption  
 117 enables a worst-case analysis of the robustness of adversarial training, which may facilitate important  
 118 theoretical insights.

### 119 3 How to Manipulate Adversarial Training

120 Previous work suggests that adversarial training could defend against both evasion attacks and delusive  
 121 attacks [39, 60]. However, in this paper, we show that adversarial training using a conventional  
 122 defense budget  $\epsilon$  may not be sufficient to provide  $\epsilon$ -robustness when confronted with stability attacks.  
 123 In particular, we present a simple theoretical model where the conventional defense scheme provably  
 124 fails when the training data is hypocritically perturbed.

125 **The binary classification task.** The data model is largely based on the setting proposed by Tsipras  
 126 et al. [63], which draws a distinction between *robust features* and *non-robust features*. Specifically, it  
 127 consists of input-label pairs  $(\mathbf{x}, y)$  sampled from a Gaussian mixture distribution  $\mathcal{D}$  as follows:

$$y \stackrel{u.a.r.}{\sim} \{-1, +1\}, \quad x_1 \sim \mathcal{N}(y, \sigma^2), \quad x_2, \dots, x_{d+1} \stackrel{i.i.d.}{\sim} \mathcal{N}(\eta y, \sigma^2), \quad (7)$$

128 where  $\eta$  is much smaller than 1 (i.e.,  $0 < \eta \ll 1$ ). Hence, samples from  $\mathcal{D}$  consist of a robust feature  
 129 ( $x_1$ ) that is *strongly* correlated with the label, and  $d$  non-robust features ( $x_2, \dots, x_{d+1}$ ) that are *very*  
 130 *weakly* correlated with it. Typically, an adversary can manipulate a large number of non-robust  
 131 features - e.g.  $d = \Theta(1/\eta^2)$  will suffice.

132 Before introducing the way to hinder robust availability, we briefly illustrate the success of adversarial  
 133 training when the training data is unperturbed.

134 **Natural and robust classifiers.** For standard classification, we consider a natural classifier:

$$f_{\text{nat}}(\mathbf{x}) := \text{sign}(\mathbf{w}_{\text{nat}}^\top \mathbf{x}), \quad \text{where } \mathbf{w}_{\text{nat}} := [1, \eta, \dots, \eta], \quad (8)$$

135 which is a minimizer of the natural risk (1) with 0-1 loss on the data (7), i.e., the Bayes optimal  
 136 classifier. However, in the adversarial setting, this natural classifier is quite brittle. Thus, it is  
 137 imperative to obtain a robust classifier:

$$f_{\text{rob}}(\mathbf{x}) := \text{sign}(\mathbf{w}_{\text{rob}}^\top \mathbf{x}), \quad \text{where } \mathbf{w}_{\text{rob}} := [1, 0, \dots, 0], \quad (9)$$

138 which relies only on the robust feature  $x_1$ .

139 **Illustration of adversarial accuracy.** In the adversarial setting, an adversary that is only allowed to  
 140 perturb each feature by a moderate  $\epsilon$  can effectively subvert the natural classifier Tsipras et al. [63]. In  
 141 particular, if  $\epsilon = 2\eta$ , an adversary can essentially force each non-robust feature to be *anti*-correlated  
 142 with the correct label. The following proposition, proved in Appendix C.1, gives the adversarial  
 143 accuracies of the natural classifier  $f_{\text{nat}}$  (8) and the robust classifier  $f_{\text{rob}}$  (9).

144 **Proposition 1.** *Let  $\epsilon = 2\eta$  and denote by  $\mathcal{A}_{\text{adv}}(f)$  the adversarial accuracy, i.e., the probability of a  
 145 classifier correctly predicting  $y$  on the data (7) under  $\ell_\infty$  perturbations. Then, we have*

$$\mathcal{A}_{\text{adv}}(f_{\text{nat}}) \leq \Pr \left\{ \mathcal{N}(0, 1) < \frac{1 - d\eta^2}{\sigma \sqrt{1 + d\eta^2}} \right\}, \quad \mathcal{A}_{\text{adv}}(f_{\text{rob}}) = \Pr \left\{ \mathcal{N}(0, 1) < \frac{1 - 2\eta}{\sigma} \right\}.$$

146 Proposition 1 implies that the adversarial accuracy of the natural classifier is  $< 50\%$  when  $d \geq 1/\eta^2$ .  
 147 Even worse, when  $\sigma \leq 1/3$  and  $d \geq 3/\eta^2$ , the adversarial accuracy of the natural classifier (8) is  
 148 always lower than 1%. In contrast, the robust classifier (9) yields a much higher adversarial accuracy  
 149 (always  $> 50\%$ ); when  $\sigma \leq (1 - 2\eta)/3$ , its adversarial accuracy will be higher than 99%.

### 150 3.1 Hypocritical Features Are Harmful

151 The results above reveal the advantages of robust classifiers over natural classifiers. Note that such  
 152 a robust classifier can be obtained by adversarial training on the original data (7). However, this  
 153 defense effect may not hold when the adversary is allowed to perturb the training data.

154 We show this by analyzing two representative perturbations: the *adversarial perturbation* [41, 60, 18]  
 155 and the *hypocritical perturbation* [59, 28, 60]. When applied into the training data, both perturbations  
 156 are effective as delusive attacks for naturally trained models. In the following, we show that the  
 157 former is harmless: adversarial training using a defense budget  $\epsilon$  on the adversarially perturbed data  
 158 can still provide test robustness. In contrast, the latter is harmful: we find that the same defense  
 159 budget can only produce non-robust classifiers when the training data is hypocritically perturbed.

160 **A harmless case.** Consider an adversary who is capable of perturbing the training data by an attack  
 161 budget  $\epsilon$ . The adversary may choose to shift each feature towards  $-y$ . Hence, the learner would see  
 162 input-label pairs  $(\mathbf{x}, y)$  sampled i.i.d. from a training distribution  $\mathcal{T}_{\text{adv}}$  as follows:

$$y \stackrel{u.a.r.}{\sim} \{-1, +1\}, \quad x_1 \sim \mathcal{N}((1 - \epsilon)y, \sigma^2), \quad x_2, \dots, x_{d+1} \stackrel{i.i.d.}{\sim} \mathcal{N}((\eta - \epsilon)y, \sigma^2), \quad (10)$$

163 where each feature of the samples from  $\mathcal{T}_{\text{adv}}$  is adversarially perturbed by a moderate  $\epsilon$ . While these  
 164 samples deviate significantly from the original distribution  $\mathcal{D}$  (7), adversarial training on them  
 165 using a defense budget  $\epsilon$  is still able to neutralize the non-robust features. Formally, in Appendix C.2  
 166 we prove the following theorem.

167 **Theorem 1** (Adversarial perturbation is harmless). *Assume that the adversarial perturbation in the*  
 168 *training data  $\mathcal{T}_{\text{adv}}$  (10) is moderate such that  $\eta/2 \leq \epsilon < 1/2$ . Then, the optimal linear  $\ell_\infty$ -robust*  
 169 *classifier obtained by minimizing the adversarial risk on  $\mathcal{T}_{\text{adv}}$  with a defense budget  $\epsilon$  is equivalent to*  
 170 *the robust classifier (9).*

171 This theorem indicates that the adversarial perturbation is harmless:  $\epsilon$ -robustness can still be obtained  
 172 by adversarial training on such perturbed training data.

173 **A harmful case.** However, this defense effect can be completely broken by the hypocritical  
 174 perturbation. That is, the adversary can instead shift each feature towards  $y$ . Hence, the learner would  
 175 see input-label pairs  $(\mathbf{x}, y)$  sampled i.i.d. from a training distribution  $\mathcal{T}_{\text{hyp}}$  as follows<sup>1</sup>:

$$y \stackrel{u.a.r.}{\sim} \{-1, +1\}, \quad x_1 \sim \mathcal{N}((1 + \epsilon)y, \sigma^2), \quad x_2, \dots, x_{d+1} \stackrel{i.i.d.}{\sim} \mathcal{N}((\eta + \epsilon)y, \sigma^2), \quad (11)$$

176 where each feature of the samples from  $\mathcal{T}_{\text{hyp}}$  is reinforced by a magnitude of  $\epsilon$ . While these samples  
 177 become more separable, adversarial training on them using the same defense budget will fail to  
 178 neutralize the hypocritically perturbed features. Consequently, the resulting classifiers will inevitably  
 179 have low adversarial accuracy. We make this formal in the following theorem proved in Appendix C.3.

180 **Theorem 2** (Hypocritical perturbation is harmful). *The optimal linear  $\ell_\infty$ -robust classifier obtained*  
 181 *by minimizing the adversarial risk on the perturbed data  $\mathcal{T}_{\text{hyp}}$  (11) with a defense budget  $\epsilon$  is equivalent*  
 182 *to the natural classifier (8).*

183 This theorem implies that the conventional defense scheme can only produce non-robust classifiers,  
 184 whose adversarial accuracy is as low as that of the natural classifier (8). That is saying, if  $\epsilon = 2\eta$ ,  
 185  $\sigma \leq 1/3$  and  $d \geq 3/\eta^2$ , the classifiers cannot get adversarial accuracy better than 1%.

186 **Implications.** As it turns out, the seemingly beneficial features in  $\mathcal{T}_{\text{hyp}}$  (11) are actually hypocritical.  
 187 Therefore, the adversary is highly motivated to hide such hypocritical features in the training data,

<sup>1</sup>To see how this relates to the hypocritical perturbation (4), let us consider the logistic loss  $\mathcal{L}(f(\mathbf{x}), y) = \log(1 + \exp(-yf(\mathbf{x})))$ , and use the natural classifier (8) as the crafting model. Then, the problem (4) has a closed-form solution  $\mathbf{p}^* = y\epsilon \cdot \text{sign}(\mathbf{w}_{\text{nat}}) = [y\epsilon, \dots, y\epsilon]$ . Applying  $\mathbf{p}^*$  to each  $\mathbf{x}$  yields the distribution  $\mathcal{T}_{\text{hyp}}$ .

188 intending to cajole an innocent learner into relying on the non-robust features. Intriguingly, we  
189 notice that the natural classifier (8) (i.e., the crafting model used to derive the distribution  $\mathcal{T}_{\text{hyp}}$ )  
190 actually has  $\eta$ -robustness. This is essentially because the non-robust features in the data (7) can resist  
191 small-magnitude perturbations by design. This motivates us to use “slightly robust” classifiers as the  
192 crafting model in practice. Indeed, our experimental results show that training the crafting model  
193 with  $0.25\epsilon$ -robustness performs the best for conducting stability attacks. This is different from the  
194 previous works [60, 18] that use naturally trained models as the crafting model for poisoning.

## 195 4 The Necessity of Large Defense Budget

196 We have shown that the hypocritical perturbation is harmful to the conventional adversarial training  
197 scheme. Fortunately, it is possible to strengthen the defense by using a larger defense budget, while  
198 the crux of the matter is how large the budget is needed.

199 We find that the minimum value of the defense budget for a successful defense depends on the specific  
200 data distribution. Let us first consider the hypocritical data in  $\mathcal{T}_{\text{hyp}}$  (11). In this case, we show that a  
201 larger defense budget is necessary in the following theorem proved in Appendix C.4.

202 **Theorem 3** ( $\epsilon + \eta$  is necessary). *The optimal linear  $\ell_\infty$ -robust classifier obtained by minimizing*  
203 *the adversarial risk on the perturbed data  $\mathcal{T}_{\text{hyp}}$  (11) with a defense budget  $\epsilon + \eta$  is equivalent to the*  
204 *robust classifier (9). Moreover, any defense budget lower than  $\epsilon + \eta$  will yield classifiers that still*  
205 *rely on all the non-robust features.*

206 This theorem implies that, in the case of the mixture Gaussian distribution under the threat of  $\epsilon$ -  
207 bounded hypocritical perturbations, the learner needs a slightly larger defense budget  $\epsilon + \eta$  to ensure  
208  $\epsilon$ -robustness.

209 While it is challenging to analyze the minimum value of the defense budget in the general case, the  
210 following theorem provides an upper bound of the budget.

211 **Theorem 4** (General case). *For any data distribution and any adversary with an attack budget  $\epsilon$ ,*  
212 *training models to minimize the adversarial risk with a defense budget  $2\epsilon$  on the perturbed data is*  
213 *sufficient to ensure  $\epsilon$ -robustness.*

214 The proof of Theorem 4 is deferred in Appendix C.5. It implies that a defense budget twice to the  
215 attack budget should be safe enough under the threat of stability attacks. Theorem 3 also suggests that  
216 the minimum budget might be much smaller than  $2\epsilon$ , and it depends on the specific attack methods  
217 and data distributions. In the following section, we empirically search for an appropriate defense  
218 budget on real-world datasets.

## 219 5 Experiments

220 In this section, we conduct comprehensive experiments to demonstrate the effectiveness of the  
221 hypocritical perturbation as stability attacks on popular benchmark datasets and the necessity of an  
222 adaptive defense for better robustness.

223 We conduct stability attacks by applying hypocritical perturbations into the training set. We focus  
224 on an  $\ell_\infty$  adversary with an *attack budget*  $\epsilon_a = 8/255$  by following [28, 74, 60, 18]. Our crafting  
225 model is adversarially trained with a *crafting budget*  $\epsilon_c = 2/255$  for 10 epochs before generating  
226 perturbations. Unless otherwise specified, we use ResNet-18 [26] as the default architecture for both  
227 the crafting model and the learning model. For adversarial training, we mainly follow the settings in  
228 previous studies [76, 65, 47]. By convention, the *defense budget* is equal to the attack budget, i.e.,  
229  $\epsilon_d = 8/255$ . More details on experimental settings are provided in Appendix D.

### 230 5.1 Benchmarking (Non-)Robustness

231 **Attack evaluation.** We compare our crafted hypocritical perturbation to existing methods, which  
232 were originally proposed as delusive attacks, including DeepConfuse (which builds an adversarial  
233 auto-encoder to generate their perturbations) [16], Unlearnable Examples (which use a min-min  
234 bi-level optimization process to pre-train their crafting model) [28], NTGA (which adopts neural  
235 tangent kernels as its crafting model) [74], and Adversarial Poisoning (whose crafting model is simply

Table 2: Test robustness (%) of PGD-AT using a defense budget  $\epsilon_d = 8/255$  on CIFAR-10.

Attack	Natural	FGSM	PGD-20	PGD-100	CW $_{\infty}$	AutoAttack
None (clean)	82.17	56.63	50.63	50.35	49.37	46.99
DeepConfuse [16]	81.25	54.14	48.25	48.02	47.34	44.79
Unlearnable Examples [28]	83.67	57.51	50.74	50.31	49.81	47.25
NTGA [74]	82.99	55.71	49.17	48.82	47.96	45.36
Adversarial Poisoning [18]	<b>77.35</b>	53.93	49.95	49.76	48.35	46.13
Hypocritical Perturbation (ours)	88.07	<b>47.93</b>	<b>37.61</b>	<b>36.96</b>	<b>38.58</b>	<b>35.44</b>

Table 3: Test robustness (%) of PGD-AT using a defense budget  $\epsilon_d = 8/255$  across different datasets.

Dataset	Attack	Natural	FGSM	PGD-20	PGD-100	CW $_{\infty}$	AutoAttack
SVHN	None	93.95	71.83	57.15	56.02	54.93	50.50
	Adv.	<b>87.50</b>	<b>56.12</b>	46.71	46.32	45.70	42.48
	Hyp.	96.06	59.41	<b>38.17</b>	<b>37.29</b>	<b>40.54</b>	<b>35.43</b>
CIFAR-100	None	56.15	31.50	28.38	28.28	26.53	24.30
	Adv.	<b>52.14</b>	28.59	26.19	26.09	24.36	22.71
	Hyp.	62.22	<b>26.38</b>	<b>21.51</b>	<b>21.13</b>	<b>21.13</b>	<b>18.74</b>
Tiny-ImageNet	None	<b>49.34</b>	25.67	22.99	22.86	20.67	18.54
	Adv.	49.52	22.93	20.01	19.91	18.75	16.83
	Hyp.	55.92	<b>20.21</b>	<b>15.61</b>	<b>15.26</b>	<b>14.99</b>	<b>12.53</b>

Table 4: Test robustness (%) of PGD-AT using a defense budget  $\epsilon_d = 8/255$  on CIFAR-10 across different architectures. Test robustness is evaluated by PGD-20. Values in parenthesis denote the accuracy on natural test data.

Attack	VGG-16	GoogLeNet	DenseNet-121	MobileNetV2	WideResNet-28-10
None	47.37 (77.15)	50.67 (83.03)	49.92 (80.08)	48.51 (80.83)	53.91 (85.81)
Adv.	44.70 (73.24)	47.72 (79.34)	48.00 (78.17)	45.90 (74.61)	51.01 (82.43)
Hyp.	<b>34.34</b> (87.20)	<b>37.03</b> (87.61)	<b>37.58</b> (88.04)	<b>35.58</b> (87.04)	<b>41.07</b> (89.14)

236 a naturally trained classifier) [18]. It is noteworthy that none of these previous works evaluated the  
 237 test robustness of their poisoned models.

238 Results using ResNet-18 on CIFAR-10 are summarized in Table 2. ‘‘Natural’’ denotes the accuracy  
 239 on natural test data. Various test-time adversarial attacks are used to evaluate test robustness,  
 240 including FGSM, PGD-20/100, CW $_{\infty}$  ( $\ell_{\infty}$  version of CW loss [9] optimized by PGD-100), and  
 241 AutoAttack (a reliable evaluation metric via an ensemble of diverse attacks [11]). We observe that  
 242 our implementation of stability attacks widely outperforms previous training-time perturbations in  
 243 degrading the test robustness of PGD-AT [39]. This demonstrates that stability attacks are indeed  
 244 harmful to the conventional defense scheme. We note that our method increases the natural accuracy.  
 245 This is reasonable, since our analysis in Section 3.1 has implied that the hypocritical perturbation can  
 246 increase model reliance on the non-robust features, which are predictive but brittle [29].

247 Moreover, we evaluate the hypocritical perturbation on other benchmark datasets including SVHN,  
 248 CIFAR-100, and Tiny-ImageNet. Both the crafting model and the victim model use the ResNet-18  
 249 architecture. Results are summarized in Table 3. ‘‘Hyp.’’ denotes the hypocritical perturbation gener-  
 250 ated by our crafting model. As a comparison, we also evaluate the adversarial perturbation gener-  
 251 ated using the same crafting model (denoted as ‘‘Adv.’’). Again, the results show that the hypocritical  
 252 perturbations are more threatening than the adversarial perturbations to standard adversarial training.  
 253 This phenomenon is consistent with our analytical results in Section 3.1.

254 Besides, we find that the hypocritical perturbation can transfer well from ResNet-18 to other architec-  
 255 tures, which reliably degrades the test robustness of a wide variety of popular architectures including  
 256 VGG-16 [54], GoogLeNet [57], DenseNet-121 [27], MobileNetV2 [49], and WideResNet-28-10 [75],  
 257 as shown in Table 4. Note that this is a completely black-box setting where the attacker has no  
 258 knowledge of the victim model’s initialization, architecture, learning rate scheduler, etc.

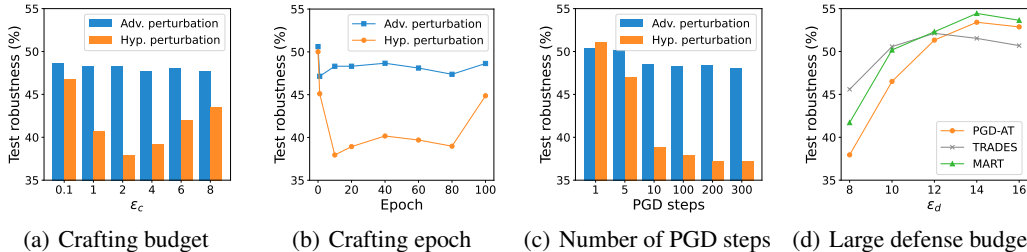


Figure 2: The ablation study experiments on CIFAR-10. Test robustness (%) is evaluated by PGD-20.

Table 5: Test robustness (%) of various adaptive defenses on the hypocritically perturbed CIFAR-10.

Defense	Natural	FGSM	PGD-20	PGD-100	$CW_\infty$	AutoAttack
PGD-AT ( $\epsilon_d = 8/255$ )	88.07	47.93	37.61	36.96	38.58	35.44
+ Random Noise	87.62	47.46	38.35	37.90	39.07	36.25
+ Gaussian Smoothing	83.95	50.96	42.80	42.34	42.41	40.07
+ Cutout	<b>88.26</b>	49.23	39.77	39.25	40.38	37.61
+ AutoAugment	86.24	48.87	40.19	39.65	37.66	35.07
PGD-AT ( $\epsilon_d = 14/255$ )	80.00	56.86	52.92	52.83	<b>50.36</b>	<b>48.63</b>
TRADES ( $\epsilon_d = 12/255$ )	79.63	55.73	51.77	51.63	48.68	47.83
MART ( $\epsilon_d = 14/255$ )	77.29	<b>57.10</b>	<b>53.82</b>	<b>53.71</b>	49.03	47.67

Table 6: Test robustness (%) of PGD-AT by adjusting the amount of clean data included in the manipulated CIFAR-10. Test robustness (%) is evaluated by PGD-20. Values in parenthesis denote the accuracy on natural test data.

Attack \ Clean proportion	0.1	0.2	0.4	0.6	0.8
None (clean subset)	30.65 (63.90)	37.99 (70.99)	44.95 (77.11)	47.17 (80.33)	49.78 (81.60)
Adversarial Perturbation	48.33 (77.71)	48.23 (76.94)	49.68 (78.54)	50.15 (82.46)	51.21 (82.03)
Hypocritical Perturbation	<b>41.51</b> (87.49)	<b>43.66</b> (88.30)	<b>46.98</b> (86.46)	<b>49.20</b> (85.29)	<b>50.56</b> (82.72)

259 **Adaptive defense.** To prevent the harm of stability attacks, our analysis in Section 4 suggests that  
 260 a larger defense budget would be helpful. We find that this is indeed the case on CIFAR-10. As  
 261 shown in Table 5, a large defense budget  $\epsilon_d = 14/255$  for PGD-AT performs significantly better  
 262 than the conventional defense budget  $\epsilon_d = 8/255$ . We also combine several data augmentations with  
 263 PGD-AT as defenses by following Fowl et al. [18]. The results show that they are beneficial, while  
 264 their improvements are inferior to PGD-AT with  $\epsilon_d = 14/255$ . In addition, we adopt other adversarial  
 265 training variants including TRADES [76] and MART [65] to defend against the hypocritical  
 266 perturbation, and find that they achieve comparable defense effects with large defense budgets.

267 Finally, we note that the adaptive defense has several limitations: *i*) robust accuracy is improved at the  
 268 cost of natural accuracy; *ii*) finding an appropriate defense budget is time-consuming for adversarial  
 269 training; *iii*) adversarial training with large budgets may lead to learning obstacles such as inherent  
 270 large sample complexity [50]. We leave the detailed study of these questions as future work.

## 271 5.2 Ablation Studies

272 In this part, we conduct a set of experiments to provide an empirical understanding of the proposed  
 273 attack. We train ResNet-18 using PGD-AT on CIFAR-10 by following the same settings described  
 274 in Appendix D unless otherwise specified.

275 **Analysis on the crafting method.** Different from previous work, we use “slightly robust” classifiers  
 276 as our crafting model. Figure 2(a) shows that this technique greatly improves the potency of the attack,  
 277 where the crafting budget  $\epsilon_c = 2/255$  performs best in degrading test robustness. We also observe  
 278 that training the crafting model for 10~80 epochs works well in Figure 2(b), and that optimizing the  
 279 crafted perturbations over 100 steps performs well in Figure 2(c). Finally, we note that Fowl et al.  
 280 [18] also tried to use adversarially trained models as the crafting model, but they failed to produce

281 effective attacks in this way. This is mainly because they adopted adversarial perturbations as poisons,  
282 which, as we observed, are inferior in degrading test robustness.

283 **Ablation on defense budget.** As discussed in Section 4, we are motivated to find the appropriate  
284 defense budget  $\epsilon_d$  in the range  $[\epsilon \sim 2\epsilon]$ . Figure 2(d) shows that the optimal defense budgets against the  
285 proposed attack for PGD-AT, TRADES, and MART are 14/255, 12/255, and 14/255, respectively.  
286 We also observe that all these adversarial training variants are inferior when using the conventional  
287 defense budget 8/255.

288 **Less data.** We follow Fowl et al. [18] to test the effectiveness of attacks by varying the proportion  
289 of clean data and perturbed data. Attacks are then considered effective if they cannot significantly  
290 increase performance over training on the clean subset alone. As shown in Table 6, the proposed  
291 attack often degrades the test robustness below what one would achieve using full clean dataset.  
292 More importantly, the hypocritical perturbations are consistently more harmful than the adversarial  
293 perturbations. This again verifies the superiority of hypocritical perturbations as stability attacks.

294 **Effect on natural training.** As a sanity check, we include the test accuracy of naturally trained  
295 models on CIFAR-10 in Appendix B Table 11. It shows that without adversarial training, the test  
296 robustness of the models becomes very poor. Thus, the goal of stability attacks is immediately  
297 achieved. On the other hand, We find that our HFs can degrade the test accuracy from 94.23%  
298 to 75.92%, though this is not the main focus of this work. We also observe that Adversarial  
299 Poisoning [18] is the most effective method in degrading the test accuracy of naturally trained models.  
300 This observation is consistent with Fowl et al. [18].

## 301 6 Related Work

302 **Adversarial training.** The presence of non-robust features has been demonstrated on popular  
303 benchmark datasets [29, 31], which naturally leads to model vulnerability to adversarial examples [63,  
304 55]. To improve test robustness against adversarial examples, adversarial training methods have been  
305 developed [23, 39, 68, 76, 61, 45, 70, 78, 58]. Usually, adversarial training using a defense budget  $\epsilon$   
306 is expected to improve model robustness against  $\epsilon$ -bounded adversarial examples. Thus, to break this  
307 defense, a direct way is to enlarge the typical  $\epsilon$ -ball used to constrain the attack; however, this may  
308 risk changing the true label [8, 62]. In this work, we aim to show that it is possible to achieve this by  
309 slightly perturbing the training data without enlarging the  $\epsilon$ -ball.

310 **Data poisoning.** Data poisoning attacks, which manipulate training data to cause the resulting  
311 models to fail during inference [3], can be divided into *availability attacks* (to degrade overall  
312 test performance) [4, 71, 40, 46, 18] and *integrity attacks* (to cause specific misclassifications) [32,  
313 10, 52, 79, 21, 51]. While the stability attacks considered in this work may be reminiscent of  
314 *backdoor attacks* [10], we note that they share several key differences. First, stability attacks  
315 aim to hinder adversarial training with well-defined  $\epsilon$ -robustness, while backdoor attacks mainly  
316 focus on embedding exploits (that can be invoked by pre-specified triggers) into naturally trained  
317 models [22, 48, 64]. Second, stability attacks only perturb the inputs slightly, while many works on  
318 backdoor attacks require mislabeling [25, 37, 44, 36]. Thus, backdoor defenses [7, 69, 35] might not  
319 be directly applied to resist stability attacks. Additional related works are discussed in Appendix A.

## 320 7 Conclusion

321 In this work, we establish a framework to study the robustness of adversarial training against stability  
322 attacks. We unveil the threat of stability attacks—small hypocritical perturbations applied into  
323 the training data suffice to hinder conventional adversarial training. The conventional defense  
324 budget  $\epsilon$  is insufficient under the threat. To resist it, we suggest that practitioners should consider a  
325 larger defense budget of no more than  $2\epsilon$ . Our theoretical analysis also explains why hypocritical  
326 perturbations work as stability attacks—they can reinforce the non-robust features and mislead  
327 the learner. Experiments demonstrate that hypocritical perturbations are harmful to conventional  
328 adversarial training on benchmark datasets, and enlarging the defense budget is essential for mitigating  
329 hypocritical perturbations. Future work includes relaxing the assumption that the adversary perturbs  
330 all the entire training set and designing more effective stability attacks against adversarial training.

## References

- 331
- 332 [1] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security:  
333 Circumventing defenses to adversarial examples. In *ICML*, 2018.
- 334 [2] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples.  
335 In *International conference on machine learning*, pages 284–293. PMLR, 2018.
- 336 [3] Battista Biggio and Fabio Roli. Wild patterns: Ten years after the rise of adversarial machine learning.  
337 *Pattern Recognition*, 84:317–331, 2018.
- 338 [4] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. In  
339 *ICML*, 2012.
- 340 [5] Battista Biggio, Iginio Corona, Davide Maiorca, Blaine Nelson, Nedim Šrndić, Pavel Laskov, Giorgio  
341 Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *ECML-PKDD*, 2013.
- 342 [6] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal,  
343 Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving  
344 cars. In *NeurIPS Deep Learning Symposium*, 2016.
- 345 [7] Eitan Borgnia, Valeriia Cherepanova, Liam Fowl, Amin Ghiasi, Jonas Geiping, Micah Goldblum, Tom  
346 Goldstein, and Arjun Gupta. Strong data augmentation sanitizes poisoning and backdoor attacks without  
347 an accuracy tradeoff. In *ICASSP*, 2021.
- 348 [8] Nicholas Carlini. A critique of the deepsec platform for security analysis of deep learning models. *arXiv*  
349 *preprint arXiv:1905.07112*, 2019.
- 350 [9] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *S&P*, 2017.
- 351 [10] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep  
352 learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- 353 [11] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of  
354 diverse parameter-free attacks. In *ICML*, 2020.
- 355 [12] Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo DeBenedetti, Nicolas Flammarion,  
356 Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness  
357 benchmark. In *NeurIPS Datasets and Benchmarks Track*, 2021. URL [https://openreview.net/  
358 forum?id=SSKZPJct7B](https://openreview.net/forum?id=SSKZPJct7B).
- 359 [13] Nilesch Dalvi, Pedro Domingos, Sumit Sanghai, and Deepak Verma. Adversarial classification. In *KDD*,  
360 2004.
- 361 [14] Edgar Dobriban, Hamed Hassani, David Hong, and Alexander Robey. Provable tradeoffs in adversarially  
362 robust classification. *arXiv preprint arXiv:2006.05161*, 2020.
- 363 [15] Ivan Evtimov, Ian Covert, Aditya Kusupati, and Tadayoshi Kohno. Disrupting model training with  
364 adversarial shortcuts. In *ICML 2021 Workshop*, 2021.
- 365 [16] Ji Feng, Qi-Zhi Cai, and Zhi-Hua Zhou. Learning to confuse: generating training time adversarial data  
366 with auto-encoder. In *NeurIPS*, 2019.
- 367 [17] Liam Fowl, Ping-yeh Chiang, Micah Goldblum, Jonas Geiping, Arpit Bansal, Wojtek Czaja, and Tom  
368 Goldstein. Preventing unauthorized use of proprietary data: Poisoning for secure dataset release. *arXiv*  
369 *preprint arXiv:2103.02683*, 2021.
- 370 [18] Liam Fowl, Micah Goldblum, Ping-yeh Chiang, Jonas Geiping, Wojtek Czaja, and Tom Goldstein.  
371 Adversarial examples make strong poisons. In *NeurIPS*, 2021.
- 372 [19] Shaopeng Fu, Fengxiang He, Yang Liu, Li Shen, and Dacheng Tao. Robust unlearnable examples: Protect-  
373 ing data privacy against adversarial learning. In *International Conference on Learning Representations*,  
374 2022. URL <https://openreview.net/forum?id=baUQQPwQiAg>.
- 375 [20] Yinghua Gao, Dongxian Wu, Jingfeng Zhang, Shu-Tao Xia, Gang Niu, and Masashi Sugiyama. Does  
376 adversarial robustness really imply backdoor vulnerability?, 2022. URL [https://openreview.net/  
377 forum?id=nG4DkcHDw\\_](https://openreview.net/forum?id=nG4DkcHDw_).
- 378 [21] Jonas Geiping, Liam H Fowl, W Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller, and Tom  
379 Goldstein. Witches’ brew: Industrial scale data poisoning via gradient matching. In *ICLR*, 2021.

- 380 [22] Micah Goldblum, Dimitris Tsipras, Chulin Xie, Xinyun Chen, Avi Schwarzschild, Dawn Song, Aleksander  
381 Madry, Bo Li, and Tom Goldstein. Dataset security for machine learning: Data poisoning, backdoor  
382 attacks, and defenses. *arXiv preprint arXiv:2012.10544*, 2020.
- 383 [23] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples.  
384 In *ICLR*, 2015.
- 385 [24] Sven Gowal, Chongli Qin, Jonathan Uesato, Timothy Mann, and Pushmeet Kohli. Uncovering the limits of  
386 adversarial training against norm-bounded adversarial examples. *arXiv preprint arXiv:2010.03593*, 2020.
- 387 [25] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine  
388 learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- 389 [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
390 In *CVPR*, 2016.
- 391 [27] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected  
392 convolutional networks. In *CVPR*, 2017.
- 393 [28] Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, and Yisen Wang. Unlearnable  
394 examples: Making personal data unexploitable. In *ICLR*, 2021.
- 395 [29] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry.  
396 Adversarial examples are not bugs, they are features. In *NeurIPS*, 2019.
- 397 [30] Adel Javanmard and Mahdi Soltanolkotabi. Precise statistical analysis of classification accuracies for  
398 adversarial training. *arXiv preprint arXiv:2010.11213*, 2020.
- 399 [31] Junho Kim, Byung-Kwan Lee, and Yong Man Ro. Distilling robust and non-robust features in adversarial  
400 examples by information bottleneck. In *NeurIPS*, 2021.
- 401 [32] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *ICML*,  
402 2017.
- 403 [33] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 404 [34] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 405 [35] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor learning:  
406 Training clean models on poisoned data. In *NeurIPS*, 2021.
- 407 [36] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor attack with  
408 sample-specific triggers. In *CVPR*, 2021.
- 409 [37] Yingqi Liu, Shiqing Ma, Youssa Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang.  
410 Trojaning attack on neural networks. 2017.
- 411 [38] Zhuoran Liu, Zhengyu Zhao, Alex Kolmus, Tijn Berns, Twan van Laarhoven, Tom Heskes, and Martha  
412 Larson. Going grayscale: The road to understanding and improving unlearnable examples. *arXiv preprint*  
413 *arXiv:2111.13244*, 2021.
- 414 [39] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards  
415 deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- 416 [40] Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee, Emil C  
417 Lupu, and Fabio Roli. Towards poisoning of deep learning algorithms with back-gradient optimization. In  
418 *ACM Workshop on Artificial Intelligence and Security*, 2017.
- 419 [41] Preetum Nakkiran. A discussion of 'adversarial examples are not bugs, they are features': Adversarial  
420 examples are just bugs, too. *Distill*, 2019. doi: 10.23915/distill.00019.5. [https://distill.pub/2019/advex-  
421 bugs-discussion/response-5](https://distill.pub/2019/advex-bugs-discussion/response-5).
- 422 [42] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits  
423 in natural images with unsupervised feature learning. 2011.
- 424 [43] James Newsome, Brad Karp, and Dawn Song. Paragraph: Thwarting signature learning by training  
425 maliciously. In *International Workshop on Recent Advances in Intrusion Detection*, 2006.
- 426 [44] Tuan Anh Nguyen and Anh Tran. Input-aware dynamic backdoor attack. In *NeurIPS*, 2020.

- 427 [45] Tianyu Pang, Xiao Yang, Yinpeng Dong, Hang Su, and Jun Zhu. Bag of tricks for adversarial training. In  
428 *ICLR*, 2021.
- 429 [46] Tianyu Pang, Xiao Yang, Yinpeng Dong, Hang Su, and Jun Zhu. Accumulative poisoning attacks on  
430 real-time data. In *NeurIPS*, 2021.
- 431 [47] Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In *ICML*, 2020.
- 432 [48] Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. Hidden trigger backdoor attacks. In  
433 *AAAI*, 2020.
- 434 [49] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2:  
435 Inverted residuals and linear bottlenecks. In *CVPR*, 2018.
- 436 [50] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially  
437 robust generalization requires more data. In *NeurIPS*, 2018.
- 438 [51] Avi Schwarzschild, Micah Goldblum, Arjun Gupta, John P Dickerson, and Tom Goldstein. Just how toxic  
439 is data poisoning? a unified benchmark for backdoor and data poisoning attacks. In *ICML*, 2021.
- 440 [52] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciuc, Christoph Studer, Tudor Dumitras, and Tom  
441 Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *NeurIPS*, 2018.
- 442 [53] Juncheng Shen, Xiaolei Zhu, and De Ma. Tensorclog: An imperceptible poisoning attack on deep neural  
443 network applications. *IEEE Access*, 2019.
- 444 [54] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recogni-  
445 tion. In *ICLR*, 2015.
- 446 [55] Jacob Springer, Melanie Mitchell, and Garrett Kenyon. A little robustness goes a long way: Leveraging  
447 robust features for targeted transfer attacks. In *NeurIPS*, 2021.
- 448 [56] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and  
449 Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2014.
- 450 [57] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru  
451 Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- 452 [58] Jihoon Tack, Sihyun Yu, Jongheon Jeong, Minseon Kim, Sung Ju Hwang, and Jinwoo Shin. Consistency  
453 regularization for adversarial robustness. In *ICML 2021 Workshop on Adversarial Machine Learning*,  
454 2021.
- 455 [59] Lue Tao, Lei Feng, Jinfeng Yi, and Songcan Chen. With false friends like these, who can notice mistakes?  
456 *arXiv preprint arXiv:2012.14738*, 2020.
- 457 [60] Lue Tao, Lei Feng, Jinfeng Yi, Sheng-Jun Huang, and Songcan Chen. Better safe than sorry: Preventing  
458 delusive adversaries with adversarial training. In *NeurIPS*, 2021.
- 459 [61] Florian Tramer and Dan Boneh. Adversarial training and robustness for multiple perturbations. In *NeurIPS*,  
460 2019.
- 461 [62] Florian Tramèr, Jens Behrmann, Nicholas Carlini, Nicolas Papernot, and Jörn-Henrik Jacobsen. Funda-  
462 mental tradeoffs between invariance and sensitivity to adversarial perturbations. In *ICML*, 2020.
- 463 [63] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robust-  
464 ness may be at odds with accuracy. In *ICLR*, 2019.
- 465 [64] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *arXiv*  
466 *preprint arXiv:1912.02771*, 2019.
- 467 [65] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving adversarial  
468 robustness requires revisiting misclassified examples. In *ICLR*, 2020.
- 469 [66] Zhirui Wang, Yifei Wang, and Yisen Wang. Fooling adversarial training with induction noise, 2022. URL  
470 <https://openreview.net/forum?id=4o1xPXaS4X>.
- 471 [67] Cheng-Hsin Weng, Yan-Ting Lee, and Shan-Hung Brandon Wu. On the trade-off between adversarial and  
472 backdoor robustness. In *NeurIPS*, 2020.

- 473 [68] Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial  
474 polytope. In *ICML*, 2018.
- 475 [69] Dongxian Wu and Yisen Wang. Adversarial neuron pruning purifies backdoored deep models. In *NeurIPS*,  
476 2021.
- 477 [70] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization.  
478 In *NeurIPS*, 2020.
- 479 [71] Huang Xiao, Battista Biggio, Gavin Brown, Giorgio Fumera, Claudia Eckert, and Fabio Roli. Is feature  
480 selection secure against training data poisoning? In *ICML*, 2015.
- 481 [72] Han Xu, Xiaorui Liu, Yaxin Li, Anil Jain, and Jiliang Tang. To be robust or to be fair: Towards fairness in  
482 adversarial training. In *ICML*, 2021.
- 483 [73] Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. Indiscriminate poisoning attacks are  
484 shortcuts. *arXiv preprint arXiv:2111.00898*, 2021.
- 485 [74] Chia-Hung Yuan and Shan-Hung Wu. Neural tangent generalization attacks. In *ICML*, 2021.
- 486 [75] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*,  
487 2016.
- 488 [76] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically  
489 principled trade-off between robustness and accuracy. In *ICML*, 2019.
- 490 [77] Jingfeng Zhang, Xilie Xu, Bo Han, Tongliang Liu, Gang Niu, Lizhen Cui, and Masashi Sugiyama. Noilin:  
491 Do noisy labels always hurt adversarial training? *arXiv preprint arXiv:2105.14676*, 2021.
- 492 [78] Jingfeng Zhang, Jianing Zhu, Gang Niu, Bo Han, Masashi Sugiyama, and Mohan Kankanhalli. Geometry-  
493 aware instance-reweighted adversarial training. In *ICLR*, 2021.
- 494 [79] Chen Zhu, W Ronny Huang, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein. Transferable  
495 clean-label poisoning attacks on deep neural nets. In *ICML*, 2019.
- 496 [80] Jianing Zhu, Jingfeng Zhang, Bo Han, Tongliang Liu, Gang Niu, Hongxia Yang, Mohan Kankanhalli, and  
497 Masashi Sugiyama. Understanding the interaction of adversarial training with noisy labels. *arXiv preprint*  
498 *arXiv:2102.03482*, 2021.

## 499 Checklist

- 500 1. For all authors...
- 501 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contribu-  
502 tions and scope? [\[Yes\]](#)
- 503 (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 2.2 and Section 7.
- 504 (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Appendix F.
- 505 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
- 506 2. If you are including theoretical results...
- 507 (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
- 508 (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#)
- 509 3. If you ran experiments...
- 510 (a) Did you include the code, data, and instructions needed to reproduce the main experimental  
511 results (either in the supplemental material or as a URL)? [\[Yes\]](#) See the supplemental material.
- 512 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?  
513 [\[Yes\]](#) See Section 5.1.
- 514 (c) Did you report error bars (e.g., with respect to the random seed after running experiments  
515 multiple times)? [\[Yes\]](#) See Appendix B.
- 516 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs,  
517 internal cluster, or cloud provider)? [\[Yes\]](#) All experiments are run on a single NVIDIA GeForce  
518 RTX 3090 GPU.
- 519 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 520 (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)

- 521 (b) Did you mention the license of the assets? [N/A]  
522 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
523 (d) Did you discuss whether and how consent was obtained from people whose data you're us-  
524 ing/curating? [N/A]  
525 (e) Did you discuss whether the data you are using/curating contains personally identifiable informa-  
526 tion or offensive content? [N/A]  
527 5. If you used crowdsourcing or conducted research with human subjects...  
528 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?  
529 [N/A]  
530 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB)  
531 approvals, if applicable? [N/A]  
532 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on  
533 participant compensation? [N/A]

## 534 A Additional Related Work

535 In this part, we discuss several independent (or concurrent) works that are related to this work.

536 Zhu et al. [80] study the effect of conventional adversarial training on differentiating noisy labels, while Zhang  
537 et al. [77] show that deliberately injected noisy labels may serve as a regularization that alleviates robust  
538 overfitting. Our results focus on the clean-label setting and provide evidence that conventional adversarial  
539 training can be hindered without modifying the labels.

540 Yu et al. [73] suggest explaining the success of availability attacks from the perspective of shortcuts. They further  
541 adopt pre-trained models to extract useful features for mitigating model reliance on the shortcuts. This direction  
542 is orthogonal to ours.

543 Liu et al. [38] improve the effectiveness of unlearnable examples [28] by generating grayscale perturbations  
544 and using data augmentations. They also conclude that conventional adversarial training will prevent a drop  
545 in accuracy measured both on clean images and adversarial images. Contrary to them, we show that, both  
546 theoretically and empirically, conventional adversarial training can be hindered by hypocritical perturbations,  
547 and we further analyze the necessity of enlarging the defense budget to resist stability attacks.

548 Gao et al. [20] revisit the trade-off between adversarial robustness and backdoor robustness [67]. They conclude  
549 that backdoor attacks are ineffective when the defense budget of adversarial training surpasses the trigger  
550 magnitude. In contrast, our results indicate that stability attacks are still harmful to adversarial training when the  
551 defense budget is not large enough. In a simple statistical setting, a defense budget  $\epsilon + \eta$  is necessary (where  $\eta$   
552 is a positive number). In the general case, a defense budget of  $2\epsilon$  is sufficient. In our experiments, a defense  
553 budget of about  $1.5\epsilon \sim 1.75\epsilon$  provides the best empirical  $\epsilon$ -robustness.

554 Wang et al. [66] argue that it is necessary to use robust features for compromising adversarial training. To this  
555 end, they adopt a relatively large attack budget  $\epsilon_a = 32/255$  for crafting their poisons (they use one type of  
556 adversarial perturbations), and show that their poisons can decrease the performance of the models trained using  
557 smaller defense budgets (such as  $\epsilon_d = 8/255$  and  $\epsilon_d = 16/255$ ). In contrast, we focus on a more realistic  
558 setting that does not require a larger attack budget. We demonstrate that it is possible to hinder adversarial  
559 training when  $\epsilon_a = \epsilon_d$ . Furthermore, we provide both theoretical and empirical results showing how to adapt the  
560 defense to maintain robustness.

561 Fu et al. [19] explore how to protect data privacy against adversarial training. The main purpose of their poisons  
562 is to compromise adversarial training by requiring the perturbation budget of their poisons to be larger than that  
563 of adversarial training. In this way, they show that the natural accuracy of the adversarially trained models can  
564 be largely decreased, let alone robust accuracy. From this perspective, our work is complementary to theirs. We  
565 pursue to not increase the attack budget of stability attacks, keeping it as small as possible. We successfully  
566 demonstrate that stability attacks are still harmful to conventional adversarial training without enlarging the  
567 attack budget. This makes the threat of stability attacks more insidious than that of Fu et al. [19].

568 On the other hand, we find that our implementation of stability attacks using hypocritical perturbations has some  
569 similarities to the robust unlearnable examples in Fu et al. [19]. Specifically, although the robust unlearnable  
570 examples are generated via a complicated min-min-max optimization process [19], we notice that their noise  
571 generator can be viewed as an adversarially trained model. This implies that the robust error-minimizing (REM)  
572 noise [19] might be useful in demonstrating the feasibility of stability attacks. To verify this, we run the source  
573 code from the authors with default hyperparameters<sup>2</sup>, and compare our crafted hypocritical perturbation with  
574 their generated noise under the setting of stability attacks. For a fair comparison, here we apply a very simple  
575 trick called EOT [2] in our method, since the trick is also used by REM [19]. The additional time cost of the  
576 EOT trick is very small and negligible.

577 Our experimental results, shown in Table 7, demonstrate that the robust error-minimizing noise is also effective  
578 as stability attacks, though it was originally proposed as a delusive attack. It is noteworthy that the robust  
579 accuracy is not evaluated in [19]. In this sense, the effectiveness of REM as an stability attack can be regarded as  
580 one of our novel findings.

581 Importantly, our method outperforms REM in terms of the robust accuracy against AutoAttack. Since AutoAttack  
582 is the most reliable evaluation metric of model robustness among the test-time attacks [11], this result indicates  
583 that our method is reliably more effective than REM in degrading model robustness.

584 It is also noteworthy that our hypocritical perturbation is significantly more efficient than the robust error-  
585 minimizing noise. For example, the time cost of our method to manipulate the CIFAR-10 dataset is 0.5 hours,  
586 whereas generating the robust error-minimizing noise for CIFAR-10 takes about 23 hours. The time cost of  
587 REM is nearly 50 times that of us!

---

<sup>2</sup><https://github.com/fshp971/robust-unlearnable-examples>

Table 7: Comparison with REM [19]: Test robustness (%) of PGD-AT using a defense budget  $\epsilon_d = 8/255$  on CIFAR-10. We report mean and standard deviation over 3 random runs.

Attack	Natural	FGSM	PGD-20	PGD-100	CW $_{\infty}$	AutoAttack
None (clean)	82.17 $\pm$ 0.71	56.63 $\pm$ 0.54	50.63 $\pm$ 0.56	50.35 $\pm$ 0.59	49.37 $\pm$ 0.57	46.99 $\pm$ 0.62
DeepConfuse [16]	81.25 $\pm$ 1.52	54.14 $\pm$ 0.63	48.25 $\pm$ 0.40	48.02 $\pm$ 0.40	47.34 $\pm$ 0.05	44.79 $\pm$ 0.36
Unlearnable Examples [28]	83.67 $\pm$ 0.86	57.51 $\pm$ 0.31	50.74 $\pm$ 0.37	50.31 $\pm$ 0.38	49.81 $\pm$ 0.24	47.25 $\pm$ 0.32
NTGA [74]	82.99 $\pm$ 0.40	55.71 $\pm$ 0.36	49.17 $\pm$ 0.27	48.82 $\pm$ 0.30	47.96 $\pm$ 0.16	45.36 $\pm$ 0.32
Adversarial Poisoning [18]	<b>77.35 <math>\pm</math> 0.43</b>	53.93 $\pm$ 0.02	49.95 $\pm$ 0.11	49.76 $\pm$ 0.08	48.35 $\pm$ 0.04	46.13 $\pm$ 0.18
REM [19]	85.63 $\pm$ 1.05	<b>42.86 <math>\pm</math> 1.09</b>	35.40 $\pm$ 0.04	35.11 $\pm$ 0.09	<b>35.24 <math>\pm</math> 0.33</b>	33.09 $\pm$ 0.24
Hypocritical Perturbation (ours)	87.60 $\pm$ 0.45	45.00 $\pm$ 0.77	<b>34.89 <math>\pm</math> 0.36</b>	<b>34.27 <math>\pm</math> 0.36</b>	36.28 $\pm$ 0.38	<b>32.79 <math>\pm</math> 0.37</b>

## 588 B Omitted Tables

Table 8: Full table of Table 2: Test robustness (%) of PGD-AT using a defense budget  $\epsilon_d = 8/255$  on CIFAR-10. We report mean and standard deviation over 3 random runs.

Attack	Natural	FGSM	PGD-20	PGD-100	CW $_{\infty}$	AutoAttack
None (clean)	82.17 $\pm$ 0.71	56.63 $\pm$ 0.54	50.63 $\pm$ 0.56	50.35 $\pm$ 0.59	49.37 $\pm$ 0.57	46.99 $\pm$ 0.62
DeepConfuse [16]	81.25 $\pm$ 1.52	54.14 $\pm$ 0.63	48.25 $\pm$ 0.40	48.02 $\pm$ 0.40	47.34 $\pm$ 0.05	44.79 $\pm$ 0.36
Unlearnable Examples [28]	83.67 $\pm$ 0.86	57.51 $\pm$ 0.31	50.74 $\pm$ 0.37	50.31 $\pm$ 0.38	49.81 $\pm$ 0.24	47.25 $\pm$ 0.32
NTGA [74]	82.99 $\pm$ 0.40	55.71 $\pm$ 0.36	49.17 $\pm$ 0.27	48.82 $\pm$ 0.30	47.96 $\pm$ 0.16	45.36 $\pm$ 0.32
Adversarial Poisoning [18]	<b>77.35 <math>\pm</math> 0.43</b>	53.93 $\pm$ 0.02	49.95 $\pm$ 0.11	49.76 $\pm$ 0.08	48.35 $\pm$ 0.04	46.13 $\pm$ 0.18
Hypocritical Perturbation (ours)	88.07 $\pm$ 1.10	<b>47.93 <math>\pm</math> 1.88</b>	<b>37.61 <math>\pm</math> 0.77</b>	<b>36.96 <math>\pm</math> 0.61</b>	<b>38.58 <math>\pm</math> 1.15</b>	<b>35.44 <math>\pm</math> 0.77</b>

Table 9: Full table of Table 3: Test robustness (%) of PGD-AT using a defense budget  $\epsilon_d = 8/255$  across different datasets. We report mean and standard deviation over 3 random runs.

Dataset	Attack	Natural	FGSM	PGD-20	PGD-100	CW $_{\infty}$	AutoAttack
SVHN	None	93.95 $\pm$ 0.21	71.83 $\pm$ 1.10	57.15 $\pm$ 0.31	56.02 $\pm$ 0.33	54.93 $\pm$ 0.19	50.50 $\pm$ 0.44
	Adv.	<b>87.50 <math>\pm</math> 0.30</b>	<b>56.12 <math>\pm</math> 0.33</b>	46.71 $\pm$ 0.25	46.32 $\pm$ 0.26	45.70 $\pm$ 0.27	42.48 $\pm$ 0.21
	Hyp.	96.06 $\pm$ 0.01	59.41 $\pm$ 0.07	<b>38.17 <math>\pm</math> 0.19</b>	<b>37.29 <math>\pm</math> 0.21</b>	<b>40.54 <math>\pm</math> 0.27</b>	<b>35.43 <math>\pm</math> 0.29</b>
CIFAR-100	None	56.15 $\pm$ 0.17	31.50 $\pm$ 0.16	28.38 $\pm$ 0.39	28.28 $\pm$ 0.40	26.53 $\pm$ 0.27	24.30 $\pm$ 0.31
	Adv.	<b>52.14 <math>\pm</math> 0.34</b>	28.59 $\pm$ 0.12	26.19 $\pm$ 0.11	26.09 $\pm$ 0.12	24.36 $\pm$ 0.09	22.71 $\pm$ 0.11
	Hyp.	62.22 $\pm$ 0.11	<b>26.38 <math>\pm</math> 0.11</b>	<b>21.51 <math>\pm</math> 0.06</b>	<b>21.13 <math>\pm</math> 0.02</b>	<b>21.13 <math>\pm</math> 0.23</b>	<b>18.74 <math>\pm</math> 0.10</b>
Tiny-ImageNet	None	<b>49.34 <math>\pm</math> 2.61</b>	25.67 $\pm$ 0.92	22.99 $\pm$ 0.37	22.86 $\pm$ 0.36	20.67 $\pm$ 0.69	18.54 $\pm$ 0.61
	Adv.	49.52 $\pm$ 0.19	22.93 $\pm$ 0.38	20.01 $\pm$ 0.24	19.91 $\pm$ 0.24	18.75 $\pm$ 0.19	16.83 $\pm$ 0.25
	Hyp.	55.92 $\pm$ 1.95	<b>20.21 <math>\pm</math> 0.84</b>	<b>15.61 <math>\pm</math> 0.31</b>	<b>15.26 <math>\pm</math> 0.26</b>	<b>14.99 <math>\pm</math> 0.73</b>	<b>12.53 <math>\pm</math> 0.57</b>

Table 10: Full table of Table 5: Test robustness (%) of various adaptive defenses on the hypocritically perturbed CIFAR-10. We report mean and standard deviation over 3 random runs.

Defense	Natural	FGSM	PGD-20	PGD-100	CW $_{\infty}$	AutoAttack
PGD-AT ( $\epsilon_d = 8/255$ )	88.07 $\pm$ 1.10	47.93 $\pm$ 1.88	37.61 $\pm$ 0.77	36.96 $\pm$ 0.61	38.58 $\pm$ 1.15	35.44 $\pm$ 0.77
+ Random Noise	87.62 $\pm$ 0.07	47.46 $\pm$ 0.08	38.35 $\pm$ 0.08	37.90 $\pm$ 0.07	39.07 $\pm$ 0.20	36.25 $\pm$ 0.14
+ Gaussian Smoothing	83.95 $\pm$ 0.27	50.96 $\pm$ 0.24	42.80 $\pm$ 0.40	42.34 $\pm$ 0.38	42.41 $\pm$ 0.19	40.07 $\pm$ 0.29
+ Cutout	<b>88.26 <math>\pm</math> 0.15</b>	49.23 $\pm$ 0.42	39.77 $\pm$ 0.26	39.25 $\pm$ 0.25	40.38 $\pm$ 0.25	37.61 $\pm$ 0.35
+ AutoAugment	86.24 $\pm$ 1.14	48.87 $\pm$ 1.01	40.19 $\pm$ 0.67	39.65 $\pm$ 0.72	37.66 $\pm$ 0.88	35.07 $\pm$ 0.88
PGD-AT ( $\epsilon_d = 14/255$ )	80.00 $\pm$ 1.91	56.86 $\pm$ 1.42	52.92 $\pm$ 0.86	52.83 $\pm$ 0.86	<b>50.36 <math>\pm</math> 1.11</b>	<b>48.63 <math>\pm</math> 0.93</b>
TRADES ( $\epsilon_d = 12/255$ )	79.63 $\pm$ 0.06	55.73 $\pm$ 0.04	51.77 $\pm$ 0.15	51.63 $\pm$ 0.15	48.68 $\pm$ 0.06	47.83 $\pm$ 0.02
MART ( $\epsilon_d = 14/255$ )	77.29 $\pm$ 0.87	<b>57.10 <math>\pm</math> 0.57</b>	<b>53.82 <math>\pm</math> 0.36</b>	<b>53.71 <math>\pm</math> 0.34</b>	49.03 $\pm$ 0.47	47.67 $\pm$ 0.51

Table 11: Test accuracy (%) of natural training on CIFAR-10. We report mean and standard deviation over 3 random runs.

Attack	Natural	PGD-20
None (clean)	94.23 $\pm$ 0.14	0.00 $\pm$ 0.00
DeepConfuse	17.22 $\pm$ 0.64	0.00 $\pm$ 0.00
Unlearnable Examples	22.72 $\pm$ 0.51	0.00 $\pm$ 0.00
NTGA	11.15 $\pm$ 0.27	0.00 $\pm$ 0.00
Adversarial Poisoning	<b>8.60 <math>\pm</math> 1.39</b>	0.00 $\pm$ 0.00
Hypocritical Perturbation	75.92 $\pm$ 1.04	0.00 $\pm$ 0.00

## 589 C Proofs

590 In this section, we provide the proofs of our theoretical results in Section 3 and Section 4.

### 591 C.1 Proof of Proposition 1

592 **Proposition 1 (restated).** *Let  $\epsilon = 2\eta$  and denote by  $\mathcal{A}_{\text{adv}}(f)$  the adversarial accuracy, i.e., the probability of a*  
 593 *classifier correctly predicting  $y$  on the data (7) under  $\ell_\infty$  perturbations. Then, we have*

$$\mathcal{A}_{\text{adv}}(f_{\text{nat}}) \leq \Pr \left\{ \mathcal{N}(0, 1) < \frac{1 - d\eta^2}{\sigma\sqrt{1 + d\eta^2}} \right\}, \quad \mathcal{A}_{\text{adv}}(f_{\text{rob}}) = \Pr \left\{ \mathcal{N}(0, 1) < \frac{1 - 2\eta}{\sigma} \right\}.$$

594

595 *Proof.* Recalling that in Equation (8), we have the natural classifier:

$$f_{\text{nat}}(\mathbf{x}) := \text{sign}(\mathbf{w}_{\text{nat}}^\top \mathbf{x}), \quad \text{where } \mathbf{w}_{\text{nat}} := [1, \eta, \dots, \eta], \quad (12)$$

596 and in Equation (9), the robust classifier is defined as:

$$f_{\text{rob}}(\mathbf{x}) := \text{sign}(\mathbf{w}_{\text{rob}}^\top \mathbf{x}), \quad \text{where } \mathbf{w}_{\text{rob}} := [1, 0, \dots, 0]. \quad (13)$$

597 Then, the adversarial accuracy of the natural classifier on the data  $\mathcal{D}$  (7) is

$$\begin{aligned} \mathcal{A}_{\text{adv}}(f_{\text{nat}}) &= 1 - \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} \left\{ \exists \|\boldsymbol{\delta}\|_\infty \leq \epsilon, f_{\text{nat}}(\mathbf{x} + \boldsymbol{\delta}) \neq y \right\} \\ &= 1 - \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [y \cdot f_{\text{nat}}(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\ &= 1 - \Pr \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} \left[ y \cdot \left( 1 \cdot (\mathcal{N}(y, \sigma^2) + \delta_1) + \sum_{i=2}^{d+1} \eta \cdot (\mathcal{N}(y\eta, \sigma^2) + \delta_i) \right) \right] < 0 \right\} \\ &\leq 1 - \Pr \left\{ y \cdot \left( 1 \cdot (\mathcal{N}(y, \sigma^2)) + \sum_{i=2}^{d+1} \eta \cdot (\mathcal{N}(y\eta, \sigma^2) - \epsilon) \right) < 0 \right\} \\ &= 1 - \Pr \left\{ \mathcal{N}(1, \sigma^2) + \eta \sum_{i=2}^{d+1} \mathcal{N}(\eta - \epsilon, \sigma^2) < 0 \right\} \\ &= \Pr \left\{ \mathcal{N}(1, \sigma^2) + \eta \sum_{i=2}^{d+1} \mathcal{N}(\eta - \epsilon, \sigma^2) > 0 \right\} \\ &= \Pr \left\{ \mathcal{N}(0, 1) < \frac{1 - d\eta^2}{\sigma\sqrt{1 + d\eta^2}} \right\}. \end{aligned} \quad (14)$$

598 Similarly, the adversarial accuracy of the robust classifier on the data  $\mathcal{D}$  (7) is

$$\begin{aligned} \mathcal{A}_{\text{adv}}(f_{\text{rob}}) &= 1 - \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} \left\{ \exists \|\boldsymbol{\delta}\|_\infty \leq \epsilon, f_{\text{rob}}(\mathbf{x} + \boldsymbol{\delta}) \neq y \right\} \\ &= 1 - \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [y \cdot f_{\text{rob}}(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\ &= 1 - \Pr \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [y \cdot (1 \cdot (\mathcal{N}(y, \sigma^2) + \delta_1))] < 0 \right\} \\ &= 1 - \Pr \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [\mathcal{N}(1, \sigma^2) + \delta_1] < 0 \right\} \\ &= 1 - \Pr \{ \mathcal{N}(1, \sigma^2) - \epsilon < 0 \} \\ &= \Pr \{ \mathcal{N}(1 - \epsilon, \sigma^2) > 0 \} \\ &= \Pr \left\{ \mathcal{N}(0, 1) < \frac{1 - 2\eta}{\sigma} \right\}. \end{aligned} \quad (15)$$

599

□

600 **C.2 Proof of Theorem 1**

601 The following theorems rely on the analytical solution of optimal linear  $\ell_\infty$ -robust classifier on mixture Gaussian  
 602 distributions. Concretely, the optimization problem is to minimize the adversarial risk on a distribution  $\widehat{\mathcal{D}}$  with a  
 603 defense budget  $\hat{\epsilon}$ :

$$\min_f \mathcal{R}_{\text{adv}}^\epsilon(f, \widehat{\mathcal{D}}), \quad \text{where } \mathcal{R}_{\text{adv}}^\epsilon(f, \widehat{\mathcal{D}}) := \mathbb{E}_{(\mathbf{x}, y) \sim \widehat{\mathcal{D}}} \left[ \max_{\|\boldsymbol{\xi}\|_\infty \leq \hat{\epsilon}} \mathbb{1} \left( \text{sign}(\mathbf{w}^\top (\mathbf{x} + \boldsymbol{\xi}) + b) \neq y \right) \right], \quad (16)$$

604 where  $f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$ , and  $\mathbb{1}(\cdot)$  denotes the indicator function.

605 We note that optimal linear robust classifiers have been obtained for certain data distributions in previous  
 606 work [63, 29, 14, 30, 72, 60]. Here, our goal is to establish similar optimal linear robust classifiers for the  
 607 classification tasks in our setting. We only employ linear classifiers, since it is highly nontrivial to consider  
 608 non-linearity for adversarial training on mixture Gaussian distributions [14].

609 **Lemma 1.** *Assume that the adversarial perturbation in data  $\mathcal{T}_{\text{adv}}$  (10) is moderate such that  $\eta/2 \leq \epsilon < 1/2$ .  
 610 Then, minimizing the adversarial risk (16) on the data  $\mathcal{T}_{\text{adv}}$  with a defense budget  $\epsilon$  can result in a classifier that  
 611 assigns 0 weight to the features  $x_i$  for  $i \geq 2$ .*

612 *Proof.* We prove the lemma by contradiction.

613 The goal is to minimize the adversarial risk on the distribution  $\mathcal{T}_{\text{adv}}$ , which can be written as follows:

$$\begin{aligned} \mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \{ \exists \|\boldsymbol{\delta}\|_\infty \leq \epsilon, f(\mathbf{x} + \boldsymbol{\delta}) \neq y \} \\ &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [y \cdot f(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\ &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \left\{ \max_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] > 0 \mid y = -1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \{y = -1\} \\ &\quad + \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] < 0 \mid y = +1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \{y = +1\} \\ &= \Pr \left\{ \underbrace{\max_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} \left[ w_1 (\mathcal{N}(\epsilon - 1, \sigma^2) + \delta_1) + \sum_{i=2}^{d+1} w_i (\mathcal{N}(\epsilon - \eta, \sigma^2) + \delta_i) + b \right]}_{\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}^{(-1)})} > 0 \right\} \cdot \frac{1}{2} \\ &\quad + \Pr \left\{ \underbrace{\min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} \left[ w_1 (\mathcal{N}(1 - \epsilon, \sigma^2) + \delta_1) + \sum_{i=2}^{d+1} w_i (\mathcal{N}(\eta - \epsilon, \sigma^2) + \delta_i) + b \right]}_{\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}^{(+1)})} < 0 \right\} \cdot \frac{1}{2} \end{aligned} \quad (17)$$

614 Consider an optimal solution  $\mathbf{w}$  in which  $w_i > 0$  for some  $i \geq 2$ . Then, we have

$$\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}^{(-1)}) = \Pr \left\{ \underbrace{\sum_{j \neq i} \max_{\|\boldsymbol{\delta}_j\| \leq \epsilon} [w_j (\mathcal{N}(\epsilon - [\mathbf{w}_{\text{nat}}]_j, \sigma^2) + \delta_j) + b]}_{\mathbb{A}} + \underbrace{\max_{\|\boldsymbol{\delta}_i\| \leq \epsilon} [w_i (\mathcal{N}(\epsilon - \eta, \sigma^2) + \delta_i)]}_{\mathbb{B}} > 0 \right\}, \quad (18)$$

615 where  $\mathbf{w}_{\text{nat}} := [1, \eta, \dots, \eta]$  as in Equation (8). Since  $w_i > 0$ ,  $\mathbb{B}$  is maximized when  $\delta_i = \epsilon$ . Thus, the  
 616 contribution of terms depending on  $w_i$  to  $\mathbb{B}$  is a normally-distributed random variable with mean  $2\epsilon - \eta$ . Since  
 617  $2\epsilon - \eta \geq 0$ , setting  $w_i$  to zero can only decrease the risk. This contradicts the optimality of  $\mathbf{w}$ . Formally,

$$\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}^{(-1)}) = \Pr \{ \mathbb{A} + w_i \mathcal{N}(2\epsilon - \eta, \sigma^2) > 0 \} > \Pr \{ \mathbb{A} > 0 \}. \quad (19)$$

618 We can also assume  $w_i < 0$  and similar contradiction holds. Therefore, minimizing the adversarial risk on  $\mathcal{T}_{\text{adv}}$   
 619 leads to  $w_i = 0$  for  $i \geq 2$ .  $\square$

620 **Lemma 2.** *Assume that the adversarial perturbation in data  $\mathcal{T}_{\text{adv}}$  (10) is moderate such that  $\eta/2 \leq \epsilon < 1/2$ .  
 621 Then, minimizing the adversarial risk (16) on the data  $\mathcal{T}_{\text{adv}}$  with a defense budget  $\epsilon$  results in a classifier that  
 622 assigns a positive weight to the feature  $x_1$ .*

623 *Proof.* We prove the lemma by contradiction.

624 The goal is to minimize the adversarial risk on the distribution  $\mathcal{T}_{\text{adv}}$ , which has been written in Equation (17).

625 Consider an optimal solution  $\mathbf{w}$  in which  $w_1 \leq 0$ . Then, we have

$$\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}^{(-1)}) = \Pr \left\{ \underbrace{\sum_{j=2}^{d+1} \max_{\|\delta_j\| \leq \epsilon} [w_j(\mathcal{N}(\epsilon - \eta, \sigma^2) + \delta_j) + b]}_{\mathbb{C}} + \underbrace{\max_{\|\delta_1\| \leq \epsilon} [w_1(\mathcal{N}(\epsilon - 1, \sigma^2) + \delta_1)]}_{\mathbb{D}} > 0 \right\}. \quad (20)$$

626 Since  $w_1 \leq 0$ ,  $\mathbb{D}$  is maximized when  $\delta_1 = -\epsilon$ . Thus, the contribution of the term depending on  $w_1$  to  $\mathbb{D}$  is a  
 627 normally-distributed random variable with mean  $-1$ . Since the mean is negative, setting  $w_1$  to be positive can  
 628 decrease the risk. This contradicts the optimality of  $\mathbf{w}$ . Formally,

$$\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}^{(-1)}) = \Pr \{ \mathbb{C} + w_1 \mathcal{N}(-\eta, \sigma^2) > 0 \} > \Pr \{ \mathbb{C} + p \mathcal{N}(-\eta, \sigma^2) > 0 \}, \quad (21)$$

629 where  $p > 0$  is any positive number. Therefore, minimizing the adversarial risk on  $\mathcal{T}_{\text{adv}}$  leads to  $w_1 > 0$ .  $\square$

630 **Theorem 1 (restated).** Assume that the adversarial perturbation in the training data  $\mathcal{T}_{\text{adv}}$  (10) is moderate such  
 631 that  $\eta/2 \leq \epsilon < 1/2$ . Then, the optimal linear  $\ell_\infty$ -robust classifier obtained by minimizing the adversarial risk  
 632 on  $\mathcal{T}_{\text{adv}}$  with a defense budget  $\epsilon$  is equivalent to the robust classifier (9).

633 *Proof.* By Lemma 1 and Lemma 2, we have  $w_1 > 0$  and  $w_i = 0$  ( $i \geq 2$ ) for an optimal linear  $\ell_\infty$ -robust  
 634 classifier. Then, the adversarial risk on the distribution  $\mathcal{T}_{\text{adv}}$  can be simplified by solving the inner maximization  
 635 problem first. Formally,

$$\begin{aligned} \mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \{ \exists \|\boldsymbol{\delta}\|_\infty \leq \epsilon, f(\mathbf{x} + \boldsymbol{\delta}) \neq y \} \\ &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [y \cdot f(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\ &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \left\{ \max_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] > 0 \mid y = -1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \{ y = -1 \} \\ &\quad + \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] < 0 \mid y = +1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{adv}}} \{ y = +1 \} \\ &= \Pr \left\{ \max_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [w_1(\mathcal{N}(\epsilon - 1, \sigma^2) + \delta_1) + b] > 0 \right\} \cdot \frac{1}{2} \\ &\quad + \Pr \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [w_1(\mathcal{N}(1 - \epsilon, \sigma^2) + \delta_1) + b] < 0 \right\} \cdot \frac{1}{2} \\ &= \Pr \{ w_1 \mathcal{N}(2\epsilon - 1, \sigma^2) + b > 0 \} \cdot \frac{1}{2} \\ &\quad + \Pr \{ w_1 \mathcal{N}(1 - 2\epsilon, \sigma^2) + b < 0 \} \cdot \frac{1}{2}, \end{aligned} \quad (22)$$

636 which is equivalent to the natural risk on a mixture Gaussian distribution  $\mathcal{D}_{\text{tmp}} : \mathbf{x} \sim \mathcal{N}(y \cdot \boldsymbol{\mu}_{\text{tmp}}, \sigma^2 \mathbf{I})$ , where  
 637  $\boldsymbol{\mu}_{\text{tmp}} = (1 - 2\epsilon, 0, \dots, 0)$ . We note that the Bayes optimal classifier for  $\mathcal{D}_{\text{tmp}}$  is  $f_{\text{tmp}}(\mathbf{x}) = \text{sign}(\boldsymbol{\mu}_{\text{tmp}}^\top \mathbf{x})$ .  
 638 Specifically, the natural risk

$$\begin{aligned} \mathcal{R}_{\text{adv}}^0(f, \mathcal{D}_{\text{tmp}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{D}_{\text{tmp}}} \{ f(\mathbf{x}) \neq y \} \\ &= \Pr_{(\mathbf{x}, y) \sim \mathcal{D}_{\text{tmp}}} \{ y \cdot f(\mathbf{x}) < 0 \} \\ &= \Pr \{ w_1 \mathcal{N}(2\epsilon - 1, \sigma^2) + b > 0 \} \cdot \frac{1}{2} \\ &\quad + \Pr \{ w_1 \mathcal{N}(1 - 2\epsilon, \sigma^2) + b < 0 \} \cdot \frac{1}{2}, \end{aligned} \quad (23)$$

639 which is minimized when  $w_1 = 1 - 2\epsilon > 0$  and  $b = 0$ . That is, minimizing the adversarial risk  $\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}})$   
 640 can lead to an optimal linear  $\ell_\infty$ -robust classifier  $f_{\text{tmp}}(\mathbf{x})$ . Meanwhile,  $f_{\text{tmp}}(\mathbf{x})$  is equivalent to the robust  
 641 classifier (9). This concludes the proof of the theorem.  $\square$

642 **C.3 Proof of Theorem 2**

643 **Lemma 3.** *Minimizing the adversarial risk (16) on the data  $\mathcal{T}_{\text{hyp}}$  (11) with a defense budget  $\epsilon$  results in a*  
 644 *classifier that assigns positive weights to the features  $x_i$  for  $i \geq 1$ .*

645 *Proof.* We prove the lemma by contradiction.

646 The goal is to minimize the adversarial risk on the distribution  $\mathcal{T}_{\text{hyp}}$ , which can be written as follows:

$$\begin{aligned}
 \mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{hyp}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{ \exists \|\boldsymbol{\delta}\|_\infty \leq \epsilon, f(\mathbf{x} + \boldsymbol{\delta}) \neq y \} \\
 &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [y \cdot f(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\
 &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \max_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] > 0 \mid y = -1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = -1\} \\
 &\quad + \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] < 0 \mid y = +1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = +1\} \\
 &= \Pr \left\{ \underbrace{\max_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} \left[ w_1(\mathcal{N}(-1 - \epsilon, \sigma^2) + \delta_1) + \sum_{i=2}^{d+1} w_i(\mathcal{N}(-\eta - \epsilon, \sigma^2) + \delta_i) + b \right]}_{\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{hyp}}^{(-1)})} > 0 \right\} \cdot \frac{1}{2} \\
 &\quad + \Pr \left\{ \underbrace{\min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} \left[ w_1(\mathcal{N}(1 + \epsilon, \sigma^2) + \delta_1) + \sum_{i=2}^{d+1} w_i(\mathcal{N}(\eta + \epsilon, \sigma^2) + \delta_i) + b \right]}_{\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{hyp}}^{(+1)})} < 0 \right\} \cdot \frac{1}{2}
 \end{aligned} \tag{24}$$

647 Consider an optimal solution  $\mathbf{w}$  in which  $w_i \leq 0$  for some  $i \geq 1$ . Then, we have

$$\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{hyp}}^{(-1)}) = \Pr \left\{ \underbrace{\sum_{j \neq i} \max_{\|\boldsymbol{\delta}_j\| \leq \epsilon} [w_j(\mathcal{N}(-[\mathbf{w}_{\text{nat}}]_j - \epsilon, \sigma^2) + \delta_j) + b]}_{\mathbb{G}} + \underbrace{\max_{\|\boldsymbol{\delta}_i\| \leq \epsilon} [w_i(\mathcal{N}(-[\mathbf{w}_{\text{nat}}]_i - \epsilon, \sigma^2) + \delta_i)]}_{\mathbb{H}} > 0 \right\}, \tag{25}$$

648 where  $\mathbf{w}_{\text{nat}} := [1, \eta, \dots, \eta]$  as in Equation (8). Since  $w_i \leq 0$ ,  $\mathbb{H}$  is maximized when  $\delta_i = -\epsilon$ . Thus, the  
 649 contribution of terms depending on  $w_i$  to  $\mathbb{H}$  is a normally-distributed random variable with mean  $-[\mathbf{w}_{\text{nat}}]_i - 2\epsilon$ .  
 650 Since the mean is negative, setting  $w_i$  to be positive can decrease the risk. This contradicts the optimality of  $\mathbf{w}$ .  
 651 Formally,

$$\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{adv}}^{(-1)}) = \Pr \{ \mathbb{G} + w_i \mathcal{N}(-[\mathbf{w}_{\text{nat}}]_i - 2\epsilon, \sigma^2) > 0 \} > \Pr \{ \mathbb{G} + p \mathcal{N}(-[\mathbf{w}_{\text{nat}}]_i - 2\epsilon, \sigma^2) > 0 \}, \tag{26}$$

652 where  $p > 0$  is any positive number. Therefore, minimizing the adversarial risk on  $\mathcal{T}_{\text{hyp}}$  leads to  $w_i > 0$  for  
 653  $i \geq 1$ .  $\square$

654 **Theorem 2 (restated).** *The optimal linear  $\ell_\infty$ -robust classifier obtained by minimizing the adversarial risk on*  
 655 *the perturbed data  $\mathcal{T}_{\text{hyp}}$  (11) with a defense budget  $\epsilon$  is equivalent to the natural classifier (8).*

656 *Proof.* By Lemma 3, we have  $w_i > 0$  for  $i \geq 1$  for an optimal linear  $\ell_\infty$ -robust classifier. Then, we have

$$\begin{aligned}
\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{hyp}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{ \exists \|\boldsymbol{\delta}\|_\infty \leq \epsilon, f(\mathbf{x} + \boldsymbol{\delta}) \neq y \} \\
&= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [y \cdot f(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\
&= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \max_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] > 0 \mid y = -1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = -1\} \\
&\quad + \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_\infty \leq \epsilon} [f(\mathbf{x} + \boldsymbol{\delta})] < 0 \mid y = +1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = +1\} \\
&= \Pr \left\{ \max_{\|\boldsymbol{\delta}_1\|_\infty \leq \epsilon} [w_1 \mathcal{N}(-1 - \epsilon, \sigma^2) + \delta_1] + \sum_{i=2}^{d+1} \max_{\|\boldsymbol{\delta}_i\|_\infty \leq \epsilon} [w_i \mathcal{N}(-\eta - \epsilon) + \delta_i] + b > 0 \right\} \cdot \frac{1}{2} \\
&\quad + \Pr \left\{ \min_{\|\boldsymbol{\delta}_1\|_\infty \leq \epsilon} [w_1 \mathcal{N}(1 + \epsilon, \sigma^2) + \delta_1] + \sum_{i=2}^{d+1} \min_{\|\boldsymbol{\delta}_i\|_\infty \leq \epsilon} [w_i \mathcal{N}(\eta + \epsilon) + \delta_i] + b < 0 \right\} \cdot \frac{1}{2} \\
&= \Pr \left\{ w_1 \mathcal{N}(-1, \sigma^2) + \sum_{i=2}^{d+1} w_i \mathcal{N}(-\eta, \sigma^2) + b > 0 \right\} \cdot \frac{1}{2} \\
&\quad + \Pr \left\{ w_1 \mathcal{N}(1, \sigma^2) + \sum_{i=2}^{d+1} w_i \mathcal{N}(\eta, \sigma^2) + b < 0 \right\} \cdot \frac{1}{2}, \tag{27}
\end{aligned}$$

657 which is equivalent to the natural risk on the mixture Gaussian distribution  $\mathcal{D} : \mathbf{x} \sim \mathcal{N}(y \cdot \mathbf{w}_{\text{nat}}, \sigma^2 \mathbf{I})$ , where  
658  $\mathbf{w}_{\text{nat}} = (1, \eta, \dots, \eta)$ . We note that the Bayes optimal classifier for  $\mathcal{D}$  is  $f_{\text{nat}}(\mathbf{x}) = \text{sign}(\mathbf{w}_{\text{nat}}^\top \mathbf{x})$ . Specifically,  
659 the natural risk

$$\begin{aligned}
\mathcal{R}_{\text{adv}}^0(f, \mathcal{D}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} \{f(\mathbf{x}) \neq y\} \\
&= \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} \{y \cdot f(\mathbf{x}) < 0\} \\
&= \Pr \left\{ w_1 \mathcal{N}(-1, \sigma^2) + \sum_{i=2}^{d+1} w_i \mathcal{N}(-\eta, \sigma^2) + b > 0 \right\} \cdot \frac{1}{2} \tag{28} \\
&\quad + \Pr \left\{ w_1 \mathcal{N}(1, \sigma^2) + \sum_{i=2}^{d+1} w_i \mathcal{N}(\eta, \sigma^2) + b < 0 \right\} \cdot \frac{1}{2},
\end{aligned}$$

660 which is minimized when  $w_1 = 1$ ,  $w_i = \eta$  for  $i \geq 2$ , and  $b = 0$ . That is, minimizing the adversarial risk  
661  $\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}_{\text{hyp}})$  can lead to an optimal linear  $\ell_\infty$ -robust classifier  $f_{\text{nat}}(\mathbf{x})$ , which is equivalent to the natural  
662 classifier (8). This concludes the proof of the theorem.

663 □

#### 664 C.4 Proof of Theorem 3

665 **Lemma 4.** *Minimizing the adversarial risk (16) on the data  $\mathcal{T}_{\text{hyp}}$  (11) with a defense budget  $\epsilon + \eta$  can result in*  
666 *a classifier that assigns 0 weight to the features  $x_i$  for  $i \geq 2$ .*

667 *Proof.* The goal is to minimize the adversarial risk on the distribution  $\mathcal{T}_{\text{hyp}}$ , which can be written as follows:

$$\begin{aligned}
\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{ \exists \|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta, f(\mathbf{x} + \boldsymbol{\delta}) \neq y \} \\
&= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [y \cdot f(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\
&= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \max_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [f(\mathbf{x} + \boldsymbol{\delta})] > 0 \mid y = -1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = -1\} \\
&\quad + \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [f(\mathbf{x} + \boldsymbol{\delta})] < 0 \mid y = +1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = +1\} \\
&= \Pr \left\{ \underbrace{\max_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} \left[ w_1 \mathcal{N}(-1 - \epsilon, \sigma^2) + \delta_1 + \sum_{i=2}^{d+1} w_i \mathcal{N}(-\eta - \epsilon, \sigma^2) + \delta_i + b \right]}_{\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}}^{(-1)})} > 0 \right\} \cdot \frac{1}{2} \\
&\quad + \Pr \left\{ \underbrace{\min_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} \left[ w_1 \mathcal{N}(1 + \epsilon, \sigma^2) + \delta_1 + \sum_{i=2}^{d+1} w_i \mathcal{N}(\eta + \epsilon, \sigma^2) + \delta_i + b \right]}_{\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}}^{(+1)})} < 0 \right\} \cdot \frac{1}{2}
\end{aligned} \tag{29}$$

668 Consider an optimal solution  $\mathbf{w}$  in which  $w_i > 0$  for some  $i \geq 2$ . Then, we have

$$\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}}^{(-1)}) = \Pr \left\{ \underbrace{\sum_{j \neq i} \max_{\|\boldsymbol{\delta}_j\| \leq \epsilon + \eta} [w_j \mathcal{N}(-[\mathbf{w}_{\text{nat}}]_j - \epsilon, \sigma^2) + \delta_j + b]}_{\mathbb{I}} + \underbrace{\max_{\|\delta_i\| \leq \epsilon + \eta} [w_i \mathcal{N}(-\eta - \epsilon, \sigma^2) + \delta_i]}_{\mathbb{J}} > 0 \right\}, \tag{30}$$

669 where  $\mathbf{w}_{\text{nat}} := [1, \eta, \dots, \eta]$ . Since  $w_i > 0$ ,  $\mathbb{J}$  is maximized when  $\delta_i = \epsilon + \eta$ . Thus, the contribution of terms  
670 depending on  $w_i$  to  $\mathbb{J}$  is a normally-distributed random variable with mean 0. Thus, setting  $w_i$  to zero will not  
671 increase the risk. Formally, we have

$$\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}}^{(-1)}) = \Pr \{ \mathbb{I} + w_i \mathcal{N}(0, \sigma^2) > 0 \} \geq \Pr \{ \mathbb{I} > 0 \}. \tag{31}$$

672 We can also assume  $w_i < 0$  and a similar argument holds. Similar arguments also hold for  $\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}}^{(+1)})$ .  
673 Therefore, minimizing the adversarial risk on  $\mathcal{T}_{\text{hyp}}$  can lead to  $w_i = 0$  for  $i \geq 2$ .  $\square$

674 **Theorem 3 (restated).** *The optimal linear  $\ell_{\infty}$ -robust classifier obtained by minimizing the adversarial risk on  
675 the perturbed data  $\mathcal{T}_{\text{hyp}}$  (11) with a defense budget  $\epsilon + \eta$  is equivalent to the robust classifier (9). Moreover, any  
676 defense budget lower than  $\epsilon + \eta$  will yield classifiers that still rely on all the non-robust features.*

677 *Proof.* By Lemma 4, we have  $w_i = 0$  ( $i \geq 2$ ) for an optimal linear  $\ell_{\infty}$ -robust classifier. Also, the robust  
678 classifier will assign a positive weight to the first feature. This is similar to the case in Lemma 2 and we omit the  
679 proof here. Then, we have

$$\begin{aligned}
\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{ \exists \|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta, f(\mathbf{x} + \boldsymbol{\delta}) \neq y \} \\
&= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [y \cdot f(\mathbf{x} + \boldsymbol{\delta})] < 0 \right\} \\
&= \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \max_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [f(\mathbf{x} + \boldsymbol{\delta})] > 0 \mid y = -1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = -1\} \\
&\quad + \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \left\{ \min_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [f(\mathbf{x} + \boldsymbol{\delta})] < 0 \mid y = +1 \right\} \cdot \Pr_{(\mathbf{x}, y) \sim \mathcal{T}_{\text{hyp}}} \{y = +1\} \\
&= \Pr \left\{ \max_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [w_1 \mathcal{N}(-1 - \epsilon, \sigma^2) + \delta_1 + b] > 0 \right\} \cdot \frac{1}{2} \\
&\quad + \Pr \left\{ \min_{\|\boldsymbol{\delta}\|_{\infty} \leq \epsilon + \eta} [w_1 \mathcal{N}(1 + \epsilon, \sigma^2) + \delta_1 + b] < 0 \right\} \cdot \frac{1}{2} \\
&= \Pr \{ w_1 \mathcal{N}(-1 - \eta, \sigma^2) + b > 0 \} \cdot \frac{1}{2} \\
&\quad + \Pr \{ w_1 \mathcal{N}(1 - \eta, \sigma^2) + b < 0 \} \cdot \frac{1}{2},
\end{aligned} \tag{32}$$

680 which is equivalent to the natural risk on a mixture Gaussian distribution  $\mathcal{D}_{\text{tmp}} : \mathbf{x} \sim \mathcal{N}(y \cdot \boldsymbol{\mu}_{\text{tmp}}, \sigma^2 \mathbf{I})$ ,  
 681 where  $\boldsymbol{\mu}_{\text{tmp}} = (1 - \eta, 0, \dots, 0)$ . We note that the Bayes optimal classifier for  $\mathcal{D}_{\text{tmp}}$  is  $f_{\text{tmp}}(\mathbf{x}) = \text{sign}(\boldsymbol{\mu}_{\text{tmp}}^\top \mathbf{x})$ .  
 682 Specifically, the natural risk

$$\begin{aligned} \mathcal{R}_{\text{adv}}^0(f, \mathcal{D}_{\text{tmp}}) &= \Pr_{(\mathbf{x}, y) \sim \mathcal{D}_{\text{tmp}}} \{f(\mathbf{x}) \neq y\} \\ &= \Pr_{(\mathbf{x}, y) \sim \mathcal{D}_{\text{tmp}}} \{y \cdot f(\mathbf{x}) < 0\} \\ &= \Pr \{w_1 \mathcal{N}(-1 - \eta, \sigma^2) + b > 0\} \cdot \frac{1}{2} \\ &\quad + \Pr \{w_1 \mathcal{N}(1 - \eta, \sigma^2) + b < 0\} \cdot \frac{1}{2}, \end{aligned} \tag{33}$$

683 which is minimized when  $w_1 = 1 - \eta > 0$  and  $b = 0$ . That is, minimizing the adversarial risk  $\mathcal{R}_{\text{adv}}^{\epsilon+\eta}(f, \mathcal{T}_{\text{hyp}})$   
 684 can lead to an optimal linear  $\ell_\infty$ -robust classifier  $f_{\text{tmp}}(\mathbf{x})$ , which is equivalent to the robust classifier (9).

685 Moreover, when the defense budget  $\epsilon_d$  is less than  $\epsilon + \eta$ , the condition in Lemma 4 no longer holds. Instead, in  
 686 this case, the robust classifier will assign positive weights to the features (i.e.,  $w_i > 0$  for  $i \geq 1$ ). This is similar  
 687 to the case in Lemma 3, and thus we omit the proof here. Consequently, this yields classifiers that still rely on all  
 688 the non-robust features. □

689

## 690 C.5 Proof of Theorem 4

691 **Theorem 4 (restated).** *For any data distribution and any adversary with an attack budget  $\epsilon$ , training models to*  
 692 *minimize the adversarial risk with a defense budget  $2\epsilon$  on the perturbed data is sufficient to ensure  $\epsilon$ -robustness.*

693 *Proof.* For clarity, we rewrite the adversarial risk in (2) with a defense budget  $\epsilon$  as follows:

$$\mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}) := \sum_{(\mathbf{x}, y) \in \mathcal{T}} \left[ \max_{\|\boldsymbol{\delta}\| \leq \epsilon} \mathcal{L}(f(\mathbf{x} + \boldsymbol{\delta}), y) \right], \tag{34}$$

694 where  $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  denotes the empirical training data.

695 Consider any adversary with an attack budget  $\epsilon$ , who can perturb  $\mathbf{x}$  to  $\mathbf{x} + \mathbf{p}$  such that  $\|\mathbf{p}\| \leq \epsilon$ . Then, the  
 696 learner will receive a perturbed version of training data  $\mathcal{T}' = \{(\mathbf{x}_i + \mathbf{p}_i, y_i)\}_{i=1}^n$ .

697 For any perturbed data point  $(\mathbf{x}_i + \mathbf{p}_i, y_i)$ , we have

$$\begin{aligned} \max_{\|\boldsymbol{\delta}\| \leq 2\epsilon} \mathcal{L}(f(\mathbf{x}_i + \mathbf{p}_i + \boldsymbol{\delta}), y_i) &= \max_{\|\boldsymbol{\delta}\| \leq \epsilon, \|\boldsymbol{\xi}\| \leq \epsilon} \mathcal{L}(f(\mathbf{x}_i + \mathbf{p}_i + \boldsymbol{\delta} + \boldsymbol{\xi}), y_i) \\ &\geq \max_{\|\boldsymbol{\delta}\| \leq \epsilon} \mathcal{L}(f(\mathbf{x}_i + \mathbf{p}_i + \boldsymbol{\delta} - \mathbf{p}_i), y_i) \\ &= \max_{\|\boldsymbol{\delta}\| \leq \epsilon} \mathcal{L}(f(\mathbf{x}_i + \boldsymbol{\delta}), y_i). \end{aligned} \tag{35}$$

698 By summarizing the training points, we have

$$\mathcal{R}_{\text{adv}}^{2\epsilon}(f, \mathcal{T}') \geq \mathcal{R}_{\text{adv}}^\epsilon(f, \mathcal{T}). \tag{36}$$

699 That is, the adversarial risk with a defense budget  $2\epsilon$  on the perturbed data is an upper bound of the adversarial  
 700 risk with a defense budget  $\epsilon$  on the original data. Therefore, a defense budget  $2\epsilon$  is sufficient to ensure the  
 701 learning of  $\epsilon$ -robustness. □

## 702 D Experimental Settings

703 **Adversary capability.** We focus on the clean-label setting, where an adversary can only provide correctly  
 704 labeled but misleading training data. In this setting, the main constraint is to craft perturbations as small as  
 705 possible [16]. Thus, we consider an  $\ell_\infty$  adversary with an *attack budget*  $\epsilon_a = 8/255$  by following Huang et al.  
 706 [28], Yuan and Wu [74], Tao et al. [60], Fowl et al. [18]. We note that this constraint is consistent with common  
 707 research on test-time adversarial examples [1].

708 **Crafting details.** We conduct stability attacks by applying the hypocritical perturbation into the training  
 709 set. Unless otherwise specified, we craft the perturbations by solving the error-minimizing objective (4) with  
 710 100 steps of PGD, where a step size of  $0.8/255$  is used by following Fowl et al. [18]. Our crafting model is  
 711 adversarially trained with a *crafting budget*  $\epsilon_c = 0.25\epsilon_a$  for 10 epochs before generating perturbations. That is,  
 712 setting  $\epsilon_c = 2/255$  performs best, as shown in Figure 2(a).

713 **Training details.** We evaluate the effectiveness of the hypocritical perturbation on benchmark datasets  
 714 including CIFAR-10/100 [33], SVHN [42], and Tiny-ImageNet [34]. Unless otherwise specified, we use ResNet-  
 715 18 [26] as the default architecture for both the crafting model and the learning model. For adversarial training,  
 716 we mainly follow the settings in previous studies [76, 65, 47]. By convention, the *defense budget* is equal to  
 717 the attack budget, i.e.,  $\epsilon_d = 8/255$ . The networks are trained for 100 epochs using SGD with momentum 0.9,  
 718 weight decay  $5 \times 10^{-4}$ , and an initial learning rate of 0.1 that is divided by 10 at the 75-th and 90-th epoch.  
 719 Early stopping is done with holding out 1000 examples from the training set. Simple data augmentations such  
 720 as random crop and horizontal flip are applied. The inner maximization problem during adversarial training is  
 721 solved by 10-steps PGD (PGD-10) with step size  $2/255$ .

## 722 E Feature-level Analysis on CIFAR-10

723 In Section 3.1, we theoretically showed that the hypocritical perturbation can cause the poisoned model to rely  
 724 more on non-robust features, thus the natural accuracy of the adversarially trained model is increased while  
 725 the robust accuracy is decreased. In this part, we aim to provide empirical evidence on the role of non-robust  
 726 features in the success of our poisoning method on a benchmark dataset. In particular, we will demonstrate that  
 727 our hypocritical perturbation successfully makes the poisoned model learn more non-robust features.

To show this, by following Section 3.2 of Ilyas et al. [29], we construct a training set where the only features  
 that are useful for classification are the non-robust features (that are extracted from the poisoned model). The  
 standard accuracy of the classifier trained on the constructed dataset can reflect how many non-robust features  
 are learned by the poisoned model (denoted as  $f$ ). To accomplish this, we modify each input-label pair  $(\mathbf{x}, y)$   
 as follows. We select a target class  $t$  uniformly at random among classes. Then, we add a small adversarial  
 perturbation to  $\mathbf{x}$  as follows:

$$\mathbf{x}_{\text{adv}} = \arg \min_{\|\mathbf{x}' - \mathbf{x}\| \leq \epsilon} \ell(f(\mathbf{x}'), t).$$

728 The resulting input-label pairs  $(\mathbf{x}_{\text{adv}}, t)$  make up the new training set. Since the resulting inputs  $\mathbf{x}_{\text{adv}}$  are nearly  
 729 indistinguishable from the originals  $\mathbf{x}$ , the label  $t$  assigned to the modified input is simply incorrect to a human  
 730 observer. Therefore, only the non-robust features in the training set are predictive, while the non-robust features  
 731 are extracted from the poisoned model.

732 We compare the model poisoned by our hypocritical perturbation with the baseline model trained on clean data.  
 733 These two models correspond to the second row and last row in Table 2, respectively. Using these two models,  
 734 we construct two datasets in the above-mentioned manner, respectively. Then, two new predictors are trained on  
 735 the two constructed datasets, respectively, and both predictors are evaluated on clean data. Training parameters  
 736 follow exactly those adopted by Ilyas et al. [29]. Our numerical results are summarized in Table 12.

Table 12: The predictive ability of the non-robust features learned by the poisoned model.

Model for constructing the training set	Standard accuracy on the original test set (%)
The baseline model	27.46
The poisoned model	<b>56.77</b>

737 As shown in Table 12, the non-robust features learned by the poisoned model are much more predictive than  
 738 the baseline. This indicates that the effect of our poisoning method on the non-robust features learned by the  
 739 poisoned model is validated empirically.

## 740 F Broader Impact

741 The attack method in this work might be used by an agent in the real world to damage the robust availability  
 742 of a machine-learning-based system. We discourage this malicious behavior by presenting the threat model of  
 743 stability attacks to the community. We further propose an adaptive defense to mitigate this issue. The adaptive  
 744 defense would help to build a more secure and robust machine learning system in the real world. At the same  
 745 time, the adaptive defense introduces an additional time cost to search for an appropriate defense budget, which  
 746 might have a negative impact on carbon emission reduction. Furthermore, society should not be overly optimistic  
 747 about AI safety, since the current studies mostly focus on perturbations bounded by simple norms (e.g.,  $\ell_\infty$  norm  
 748 in this paper). There might exist perturbations beyond the  $\ell_p$  ball in the real world, and we are still far from  
 749 complete model robustness.