

# How Robust are Randomized Smoothing based Defenses to Data Poisoning?

Akshay Mehra<sup>1</sup>, Bhavya Kailkhura<sup>2</sup>, Pin-Yu Chen<sup>3</sup> and Jihun Hamm<sup>1</sup>

<sup>1</sup>Tulane University <sup>2</sup>Lawrence Livermore National Laboratory <sup>3</sup> IBM Research

{amehra, jhamm3}@tulane.edu, kailkhura1@llnl.gov, pin-yu.chen@ibm.com

## Abstract

*Predictions of certifiably robust classifiers remain constant in a neighborhood of a point, making them resilient to test-time attacks with a guarantee. In this work, we present a previously unrecognized threat to robust machine learning models that highlights the importance of training-data quality in achieving high certified adversarial robustness. Specifically, we propose a novel bilevel optimization based data poisoning attack that degrades the robustness guarantees of certifiably robust classifiers. Unlike other poisoning attacks that reduce the accuracy of the poisoned models on a small set of target points, our attack reduces the average certified radius (ACR) of an entire target class in the dataset. Moreover, our attack is effective even when the victim trains the models from scratch using state-of-the-art robust training methods such as Gaussian data augmentation[8], MACER[36], and SmoothAdv[29] that achieve high certified adversarial robustness. To make the attack harder to detect, we use clean-label poisoning points with imperceptible distortions. The effectiveness of the proposed method is evaluated by poisoning MNIST and CIFAR10 datasets and training deep neural networks using previously mentioned training methods and certifying the robustness with randomized smoothing. The ACR of the target class, for models trained on generated poison data, can be reduced by more than 30%. Moreover, the poisoned data is transferable to models trained with different training methods and models with different architectures.*

## 1. Introduction

Data poisoning [3, 17, 31, 32, 37] is a training-time attack where the attacker is assumed to have access to the training data on which the victim will train the model. The attacker can modify the training data in a manner that the model trained on this poisoned data performs as the attacker desires. The data hungry nature of modern machine learning methods make them vulnerable to poisoning attacks. Attackers can place the poisoned data online and wait for it to be scraped by victims trying to increase the size for their training sets. Another easy target for data poisoning is data collection by crowd sourcing where malicious users can corrupt the data

they contribute. In most cases, an attacker can modify only certain parts of the training data such as change the features or labels for a specific class or modify a small subset of the data from all classes. In this work, we assume the attacker wants to affect the performance of the victim’s models on a target class and modifies only the features of the points belonging to that class (without affecting the labels). To evade detection, the attacker is constrained to only add imperceptibly small perturbations to the points of the target class. Many previous works [26, 31, 16, 19, 37, 7, 18, 33] have shown the effectiveness of poisoning in affecting the accuracy of models trained on poisoned data compared to the accuracy achievable by training with clean data. In most works, the victim is assumed to use standard training by minimizing the empirical loss on the poisoned data to train the models and thus the attack is optimized to hurt the accuracy of standard training. However, recent research on test-time evasion attacks [5, 1, 34, 4] suggests that models trained with standard training are not robust to adversarial examples, making the assumption of victim relying on standard training to train the models for deployment questionable.

Thus, in a realistic scenario, where the aim of the victim is to deploy the model, it’s better to assume that the victim will rely on training procedures that yield classifiers which are provably robust to test-time attacks. Several recent works have proposed methods for training certifiably robust models whose predictions are guaranteed to be constant in a neighbourhood of a point. However, many of these methods [28, 13, 15, 35] do not scale to deep neural networks or large datasets, due to their high complexity. Moreover, the effect of training data quality on the performance of these certified defenses at test time remains largely unexplored. Recently, randomized smoothing (RS) based certification methods [20, 21, 8] were shown to be scalable to deep neural networks and high dimensional datasets enabling researchers to propose training procedures [29, 36] that lead to models with high certified robustness. Thus, we assume that a victim will rely on RS based certification methods to measure the certified robustness and use RS based training procedures to train the models. The fact that a victim can train with a method that improves certified adversarial robustness is

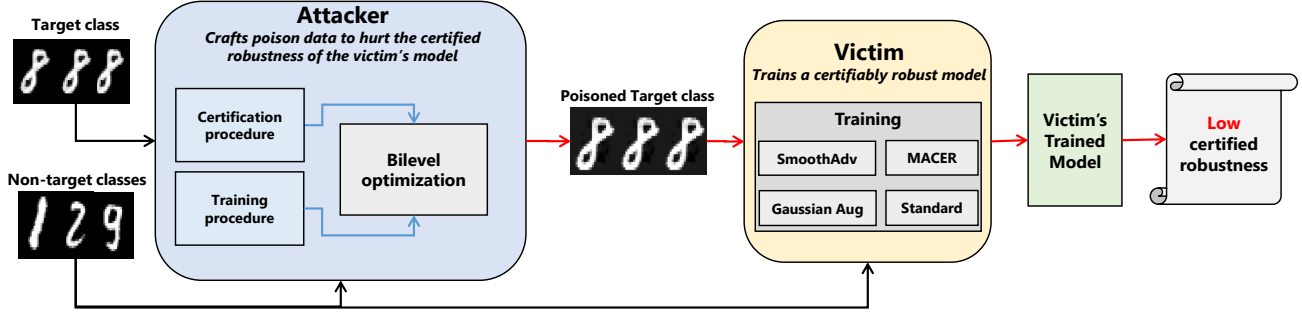


Figure 1. Overview of our poisoning against certified defenses (PACD) attack which generates poisoned data to reduce the certified robustness of the victim’s model trained with methods such as Gaussian data augmentation(GA)[8], SmoothAdv[29] and MACER[36] on a target class.

Table 1. Failure of traditional data poisoning attacks optimized against standard training in affecting the test accuracy (of target class) of models trained with certifiably robust training procedures. Details of the experiment are present in Appendix D.1. Certifiably robust training methods [8, 29, 36] are trained with  $\sigma = 0.25$  and accuracy of their base classifiers are reported.

	Training method	Model trained on clean data	Model trained on poison data
MNIST	Standard	99.28±0.01	60.08±12.6
	GA[8]	98.99±0.14	98.31±1.65
	SmoothAdv [29]	99.18±0.23	99.31±0.29
	MACER [36]	99.21±0.56	98.31±0.58
CIFAR10	Standard	92.71±1.31	0.36±0.37
	GA [8]	88.84±2.39	88.38±2.13
	SmoothAdv [29]	79.48±2.69	74.95±3.45
	MACER [36]	87.12±1.17	88.54±4.52

an immediate challenge for current poisoning attacks which optimize the poison data to affect the accuracy of models trained with standard training. Table 1 shows that poisons optimized against standard training can significantly reduce the accuracy of the victim’s model (left to right) when the victim also uses standard training (1st and 5th row). However, this poison data is rendered ineffective when the victim uses a certifiably robust training method such as [8, 29, 36].

*Are certified defenses robust to data poisoning?* We study this question and demonstrate that data poisoning is a serious concern even for certified defenses. We propose a novel data poisoning attack that can significantly compromise the certified robustness guarantees achievable from training with robust training procedures. We formulate the Poisoning Against Certified Defenses (PACD) attack as a constrained bilevel optimization problem and theoretically analyze its solution for the case when the victim uses linear classifiers. Our theoretical analysis and empirical results suggests that the decision boundary of the smoothed classifiers (used for RS) learned from the poisoned data is significantly different from the one learned from clean data there by causing a

reduction in certified radius. Our bilevel optimization based attack formulation is general since it can generate poisoned data against a model trained with any certifiably robust training method (lower-level problem) and certified with any certification procedure (upper-level problem). Fig. 1 shows the overview of the proposed PACD attack.

Unlike previous poisoning attacks that aim to reduce the accuracy of the models on a small subset of data, our attack can reduce the certified radius of an entire target class. The poison points generated by our attack have clean labels and imperceptible distortion making them difficult to detect. The poison data remains effective when the victim trains the models from scratch or uses data augmentation or weight regularization during training. Moreover, the attack points generated against a certified defense are transferable to models trained with other RS based certified defenses and to models with different architectures. This highlights the importance of training-data quality and curation for obtaining meaningful gains from certified defenses at test time, a factor not considered by current certified defense research.

Our main contributions are as follows

- We study the problem of using data poisoning attacks to affect the robustness guarantees of classifiers trained using certified defense methods. To the best of our knowledge, this is the first clean label poisoning attack that significantly reduces the certified robustness guarantees of the models trained on the poisoned dataset.
- We propose a bilevel optimization based attack which can generate poison data against several robust training and certification methods. We specifically use the attack to highlight the vulnerability of randomized smoothing based certified defenses to data poisoning.
- We demonstrate the effectiveness of our attack in reducing the certifiable robustness obtained using randomized smoothing on models trained with state-of-the-art certified defenses [8, 29, 36]. Our attack reduces the ACR of the target class by more than 30%.

## 2. Background and related work

**Randomized smoothing:** The RS procedure [8] uses a smoothed version of the original classifier  $f : \mathbb{R}^d \rightarrow \mathcal{Y}$  and certifies the adversarial robustness of the new classifier. The smoothed classifier,  $g(x) = \arg \max_c \mathbb{P}_{\eta \sim \mathcal{N}(0, \sigma^2 I)}(f(x + \eta) = c)$ , assigns  $x$  the class whose decision region  $\{x' \in \mathbb{R}^d : f(x') = c\}$  has the largest measure under the distribution  $\mathcal{N}(x, \sigma^2 I)$ , where  $\sigma$  is used for smoothing. Suppose that while classifying a point  $\mathcal{N}(x, \sigma^2 I)$ , the original classifier  $f$  returns the class  $c_A$  with probability  $p_A = \mathbb{P}(f(x + \eta) = c_A)$ , and the “runner-up” class  $c_B$  is returned with probability  $p_B = \max_{c \neq c_A} \mathbb{P}(f(x + \eta) = c)$ , then the prediction of the point  $x$  under the smoothed classifier  $g$  is robust within the radius  $r(g; \sigma) = \frac{\sigma}{2}(\Phi^{-1}(p_A) - \Phi^{-1}(p_B))$ , where  $\Phi^{-1}$  is the inverse CDF of the standard Normal distribution. In practice, Monte Carlo sampling is used to estimate a lower bound on  $p_A$  and an upper bound on  $p_B$  as its difficult to estimate the actual values for  $p_A$  and  $p_B$ . Since standard training of the base classifier does not achieve high robustness guarantees, [8] proposed to use GA based training in which the base classifier is trained on Gaussian noise corruptions of the clean data. Recent works [36, 29] showed that the certified robustness guarantees of RS can be boosted by using different training procedures. In particular, [29] proposed to train the base classifier using adversarial training where the adversarial examples are generated against the smoothed classifier. Although effective at increasing the certified radius, the method can be slow to train due to the requirement of generating adversarial examples against the smoothed classifier at every step. Another recent work [36] proposed a different training procedure which is significantly faster to train and relies on directly maximizing the certified radius for achieving high robustness guarantees. Due to their effectiveness in improving the certified robustness guarantees of machine learning models, we craft poison data against these methods. A recent attack method [11] showed that it is possible to fool a robust classifier to mislabel an input and give an incorrect certificate using perturbation large in  $\ell_p$  norm at test-time. Our work is different since we focus on train-time attacks against certified defenses using imperceptibly small perturbations to the poison data.

**Bilevel optimization:** A bilevel optimization problem has the form  $\min_{u \in \mathcal{U}} \xi(u, v^*)$  s.t.  $v^* = \arg \min_{v \in \mathcal{V}(u)} \zeta(u, v)$ , where the upper-level problem is a minimization problem with  $v$  constrained to be the optimal solution to the lower-level problem (see [2]). Our data poisoning attack is a constrained bilevel optimization problem. Although general bilevel problems are difficult to solve, under some simplifying assumptions their solution can be obtained using gradient based methods. Several methods for solving bilevel problems in machine learning have been proposed previously [9, 27, 10, 22, 30, 24] (See Appendix B for an overview). We use the method based on approximating

the hypergradient by approximately solving a linear system (ApproxGrad Alg. 2 in Appendix B) in this work. Previous works [25, 26, 24, 16, 6] have shown the effectiveness of solving bilevel optimization problem for data poisoning to affect the accuracy of models trained with standard training. Our work on the other hand proposes a bilevel optimization based formulation to generate a data poisoning attack against RS based certified defenses and shows its effectiveness against state-of-the-art robust training methods.

## 3. Poisoning against certified defenses

Here we present the bilevel formulation of our PACD attack for generating poisoned data to compromise the certified robustness guarantees of the models trained using certified defenses. Specifically, we discuss how to generate poison data against GA [8], SmoothAdv [29] and MACER [36] and affect the certified robustness guarantees obtained using RS.

### 3.1. General attack formulation

Let  $\mathcal{D}^{\text{clean}} = \{(x_i^{\text{clean}}, y_i^{\text{clean}})\}_{i=1}^{N^{\text{clean}}}$  be the clean, unalterable portion of the training set. Let  $u = \{u_1, \dots, u_n\}$  denote the attacker’s poisoning data which is added to the clean data. For clean-label attack, we require that each poison example  $u_i$  has a limited perturbation, for example,  $\|u_i - x_i^{\text{base}}\| = \|\delta_i\| \leq \epsilon$  from the base data  $x_i^{\text{base}}$  and has the same label  $y_i^{\text{base}}$ , for  $i = 1, \dots, n$ . Thus  $\mathcal{D}^{\text{poison}} = \{(u_i, y_i^{\text{base}})\}_{i=1}^{N^{\text{poison}}}$ . The goal of the attacker is to find  $u$  such that when the victim uses  $\mathcal{D}^{\text{clean}} \cup \mathcal{D}^{\text{poison}}$  to train a classifier, the certified robustness guarantees of the model on the target class ( $\mathcal{D}^{\text{val}} = \{(x_i^{\text{val}}, y_i^{\text{val}})\}_{i=1}^{N^{\text{val}}}$ ) are significantly diminished compared to a classifier trained on clean data. The attack can be formulated as follows:

$$\begin{aligned} \min_{u \in \mathcal{U}} \mathcal{R}(\mathcal{D}^{\text{val}}; \theta^*) \\ \text{s.t. } \theta^* = \arg \min_{\theta} \mathcal{L}_{\text{robust}}(\mathcal{D}^{\text{clean}} \cup \mathcal{D}^{\text{poison}}; \theta). \end{aligned} \quad (1)$$

The upper-level cost  $\mathcal{R}$  denotes a certified robustness metric such as the certified radius from RS. The goal of the upper-level problem is to compromise the certified robustness guarantees of the model trained on validation data  $\mathcal{D}^{\text{val}}$ . The solution to the lower-level problem  $\theta^*$  are the parameters of the machine learning model learnt from  $\mathcal{D}^{\text{clean}} \cup \mathcal{D}^{\text{poison}}$  using a robust training method with loss function  $\mathcal{L}_{\text{robust}}$ . The fact that any training method that achieves high certified robustness to test-time attacks and any certification procedure can be incorporated into this formulation by changing the lower- and upper-level problems, respectively, makes the attack formulation broadly applicable. Recent works [8, 29, 36] have shown RS based methods to be effective at certifying and producing robust classifiers. The scalability of these methods to large datasets and deep models make them useful for real-world application. Thus, we focus on using our poisoning attack against these methods.

### 3.2. Poison randomized smoothing based defenses

For an input at test time, RS produces a prediction from the smoothed classifier  $g$  and a radius in which this prediction remains constant. Since the certified radius of a “hard” smooth classifier  $g$  is non-differentiable, it cannot be directly incorporated in the upper-level of the attack formulation Eq. (1). To overcome this challenge, we use the “soft” smooth classifier  $\tilde{g}$  as an approximation. Similar technique has been used in [29, 36]. Let  $z_\theta : X \rightarrow \mathcal{P}(K)$  be a classifier whose last layer is softmax with parameters  $\theta$  and  $\sigma > 0$  is the noise used for smoothing, then soft smoothed classifier  $\tilde{g}_\theta$  of  $z_\theta$  is  $\tilde{g}_\theta(x) = \arg \max_{c \in \mathcal{Y}} \mathbb{E}_{\eta \sim \mathcal{N}(0, \sigma^2 I)} [z_\theta^c(x + \eta)]$ . It was shown in [36] that if the ground truth of an input  $x$  is  $y$  and  $\tilde{g}_\theta$  classifies  $x$  correctly then  $\tilde{g}_\theta$  is provably robust at  $x$ , with the certified radius  $\tilde{r}(\tilde{g}_\theta; x, y, \sigma) = \frac{\sigma}{2} [\Phi^{-1}(\mathbb{E}_\eta [z_\theta^y(x + \eta)]) - \Phi^{-1}(\max_{y' \neq y} \mathbb{E}_\eta [z_\theta^{y'}(x + \eta)])]$ . Assuming  $\tilde{r}(\tilde{g}_\theta; x, y, \sigma) = 0$  when  $x$  is misclassified, the ACR is  $R(\tilde{g}_\theta; \mathcal{D}, \sigma) = \frac{1}{|\mathcal{D}|} \sum_{(x, y) \in \mathcal{D}} \tilde{r}(\tilde{g}_\theta; x, y, \sigma)$ . Since  $\tilde{R}$  is differentiable we can use it in the upper-level of Eq. (1). The lower-level problem can be any robust training procedure and we focus on using [8, 36, 29] in this work.

**Poisoning against GA [8].** We start by showing how to generate poison data against models trained with GA [8], which was shown to yield higher certified robustness compared to models trained with standard training. In this method the classifier  $f_\theta$  is obtained by optimizing the loss function  $\mathcal{L}_{\text{GaussAug}}(\mathcal{D}; \theta, \sigma) = \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} l_{ce}(x_i + \eta, y_i; \theta)$ , where  $l_{ce}$  is the cross entropy loss and  $\eta \sim \mathcal{N}(0, \sigma^2 I)$ . To control the perturbation added to the poison data we used  $\ell_\infty$ -norm here but other norms can also be used. The bilevel formulation to generate poison data to reduce the certified robustness guarantees obtained using RS for a classifier trained with GA is as follows.

$$\begin{aligned} \min_u \quad & \tilde{R}(\tilde{g}_{\theta^*}; \mathcal{D}^{\text{val}}, \sigma) \\ \text{s.t.} \quad & \|\delta_i\|_\infty \leq \epsilon, \quad i = 1, \dots, n, \quad \text{and} \end{aligned} \quad (2)$$

$$\theta^* = \arg \min_\theta \mathcal{L}_{\text{GaussAug}}(\mathcal{D}^{\text{clean}} \cup \mathcal{D}^{\text{poison}}; \theta, \sigma).$$

**Poisoning against MACER [36].** Another recent work proposed a method for robust training by maximizing the certified radius (MACER). Their approach uses a loss function which is a combination of the classification loss and the robustness loss of the soft smoothed classifier  $\tilde{g}_\theta$ . In particular, the loss of the smoothed classifier on a point  $(x, y)$  is given by  $l_{\text{macer}}(\tilde{g}_\theta; x, y) = -\log \hat{z}_\theta^y(x) + \frac{\lambda \sigma}{2} \max\{\gamma - \tilde{\xi}_\theta(x, y), 0\} \cdot \mathbf{1}_{\tilde{g}_\theta(x) = y}$ . where  $\eta_1, \dots, \eta_k$  are  $k$  i.i.d. samples from  $\mathcal{N}(0, \sigma^2 \mathbf{I})$ ,  $\hat{z}_\theta^y(x) = \frac{1}{k} \sum_{j=1}^k z_\theta^y(x + \eta_j)$  is the empirical expectation of  $z_\theta(x + \eta)$ ,  $\tilde{\xi}_\theta(x, y) = \Phi^{-1}(\hat{z}_\theta^y(x)) - \Phi^{-1}(\max_{y' \neq y} \hat{z}_\theta^{y'}(x))$ ,  $\gamma$  is the hinge factor, and  $\lambda$  balances the accuracy and robustness trade-off. Using this we can define  $\mathcal{L}_{\text{macer}}(\mathcal{D}; \theta, \sigma) = \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} l_{\text{macer}}(\tilde{g}_\theta; x_i, y_i)$ .

---

#### Algorithm 1 Poisoning GA based certified defense [8]

---

**Input:**  $\mathcal{D}^{\text{clean}}, \mathcal{D}^{\text{base}}, \mathcal{D}^{\text{val}}$ , perturbation strength  $\epsilon$ , noise level  $\sigma$ , number of noise samples  $k$ , inverse temperature  $\alpha$ , total epochs  $P$ , lower – level epochs  $T_1$ , epochs for linear system  $T_2$

**Output:**  $\mathcal{D}^{\text{poison}}$

$\mathcal{D}^{\text{poison}} := \mathcal{D}^{\text{base}}$

**for**  $p = 0, \dots, P-1$  **do**

    Sample a mini-batch  $(x^{\text{clean}}, y^{\text{clean}}) \sim \mathcal{D}^{\text{clean}}$

    Sample a mini-batch of  $n$  points  $(x^{\text{val}}, y^{\text{val}}) \sim \mathcal{D}^{\text{val}}$

    Sample a mini-batch  $(x^{\text{poison}}, y^{\text{poison}}) \sim \mathcal{D}^{\text{poison}}$

    Pick the base samples for poison data  $(x^{\text{base}}, y^{\text{base}})$

    For each  $x_i^{\text{val}}$ , sample  $k$  i.i.d. Gaussian samples

$x_{i_1}^{\text{val}}, \dots, x_{i_n}^{\text{val}} \sim \mathcal{N}(x_i^{\text{val}}, \sigma^2 I)$

    Compute  $\tilde{z}_\theta(x_i^{\text{val}}) \leftarrow \frac{1}{k} \sum_{j=1}^k \alpha z_\theta(x_{i_j}^{\text{val}})$  for each  $i$

$\mathcal{G}_\theta := \{(x_i^{\text{val}}, y_i^{\text{val}}) : y_i^{\text{val}} = \arg \max_{c \in \mathcal{Y}} \tilde{z}_\theta^c(x_i^{\text{val}})\}$

    For each  $(x_i, y_i) \in \mathcal{G}_\theta$ , compute  $\tilde{y}_i$

$\tilde{y}_i \leftarrow \arg \max_{c \in \mathcal{Y} \setminus \{y_i\}} \tilde{z}_\theta^c(x_i)$

    For each  $(x_i, y_i) \in \mathcal{G}_\theta$ , compute  $\tilde{r}(x_i, y_i)$

$\tilde{r}(x_i, y_i) = \frac{\sigma}{2} (\Phi^{-1}(\tilde{z}_\theta^{y_i}(x_i)) - \Phi^{-1}(\tilde{z}_\theta^{\tilde{y}_i}(x_i)))$

$\xi := \frac{1}{n} \sum_{(x_i, y_i) \in \mathcal{G}_\theta} \tilde{r}(x_i, y_i)$

$\zeta := \mathcal{L}_{\text{GaussAug}}((x^{\text{clean}}, y^{\text{clean}}) \cup (x^{\text{poison}}, y^{\text{poison}}), \sigma)$

$(x^{\text{poison}}, y^{\text{poison}}) := \text{ApproxGrad}(\xi, \zeta, 1, T_1, T_2, \epsilon, x^{\text{base}})$

    Update  $\mathcal{D}^{\text{poison}}$  with  $(x^{\text{poison}}, y^{\text{poison}})$

**end for**

---

To generate poison data that reduces the robustness guarantees of classifier trained with MACER we can use the loss  $\mathcal{L}_{\text{macer}}(\mathcal{D}; \theta, \sigma)$  in the lower-level problem in Eq. (2).

**Poisoning against SmoothAdv [29].** It was shown that the certified robustness guarantees obtained from RS can be improved by training the classifiers using adversarial training with adversarial examples generated against the smooth classifier. In particular the classifier trained with SmoothAdv optimizes the following objective for a point  $(x, y)$ .  $\min_\theta \max_{\|x' - x\|_2 \leq \alpha} -\log \frac{1}{k} \sum_{j=1}^k z_\theta^y(x' + \eta_j)$  where  $\eta_1, \dots, \eta_k$  are  $k$  i.i.d. samples from  $\mathcal{N}(0, \sigma^2 \mathbf{I})$  and  $\alpha$  is the permissible  $\ell_2$  distortion to  $x$ . To generate poisoning data against SmoothAdv we must use this objective as the lower-level problem in Eq. (2). To make it easier for bilevel solvers to solve this problem we use an approximation to the min-max problem. For doing that we first compute the adversarial example  $x' = \arg \max_{\|x' - x\|_2 \leq \alpha} -\log \frac{1}{k} \sum_{j=1}^k z_\theta^y(x' + \eta_j)$  using PGD attack on the points in  $\mathcal{D}^{\text{clean}} \cup \mathcal{D}^{\text{poison}}$  and then use these examples as our new dataset to train the model parameters in the lower-level as in Eq. (2). Specifically, the lower-level problem in Eq. (2) becomes

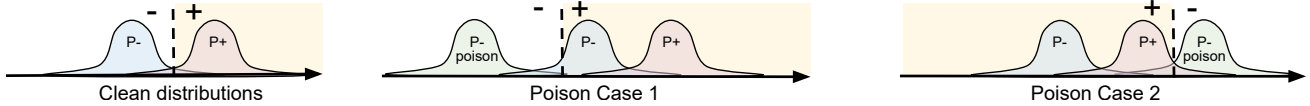


Figure 2. Analytical solutions of problem (3) with linear classifiers. The poison distribution ( $P_{\text{poison}}^-$ ) can change the decision boundary (broken line) and reduce the ACR of the clean distribution ( $P^-$ ) in two ways (Cases 1 and 2). Perturbation is exaggerated for illustration.

$\arg \min_{\theta} \mathcal{L}_{\text{GaussAug}}(\mathcal{D}_{\text{adv}}^{\text{clean}} \cup \mathcal{D}_{\text{adv}}^{\text{poison}}; \theta, \sigma)$  where  $\mathcal{D}_{\text{adv}}$  denotes the adversarial examples generated against  $\tilde{g}_{\theta}$ . We update  $\mathcal{D}_{\text{adv}}$  in every step of the bilevel optimization.

### 3.3. Generation and evaluation of poisoning attack

In this work, we focus on creating a poisoned set to compromise the certified adversarial robustness guarantees of all points in a target class. We initialize the poison data with clean data from the target class (i.e., base data) and optimize the perturbation to be added to each point by solving the bilevel problem in Eq. (2) for attack against GA based training. We use a small value of  $\epsilon$  to ensure the perturbations added are imperceptible and the poison points have clean labels when inspected visually (See Fig. 5 in the Appendix). The bilevel optimization is solved using the ApproxGrad algorithm (Alg. 2 in Appendix B). The full attack algorithm for generating poison data against GA [8] is shown in Alg. 1. Attack against other methods are generated similarly by replacing the lower-level objective ( $\zeta$  in Alg. 1) with the appropriate loss function for MACER [36] and SmoothAdv [29]. We evaluate the effect of poisoning, by training the models from scratch using GA, MACER and SmoothAdv on their respective poisoned sets and report ACR and approximate certified accuracy (points with certified  $\ell_2$  radius greater than zero) on the clean test points from the target class. Previous works [16, 31] had shown the effectiveness of poisoning by lowering the accuracy on specific target points from the test set. Our attack is also effective, under a similar setting, at reducing the certified radius for target points (Appendix D.6).

### 3.4. Analysis of poisoning with linear classifiers

To gain a deeper insight into the effect of poisoning, we analyze the analytical solution of our bilevel problem for the case of linear classifiers trained with GA. Suppose we have a one-dimensional two-class problem and the attacker’s goal is to poison the distribution of the *negative* class  $P^-$  so that the ACR ( $\bar{R}$ ) of the poisoned model on the test points of the *negative* class is reduced. Let  $\epsilon$  be the maximum permissible perturbation that can be added by the attacker to the points of the class  $P^-$ . We do not assume any specific distributions for  $P^+$  and  $P^-$  here, but only that  $\sum_i x_i^- < \sum_i x_i^+$  without loss of generality. Here  $x_i^+$  and  $x_i^-$  refer to the training points of the positive and the negative class, respectively. A linear classifier in one-dimension is either  $f(x) = 1$  iff  $x \geq t$  or  $f(x) = -1$  iff  $x \leq t$  parameterized by the threshold  $t$ . For

linear classifiers, the smoothed classifier  $g$  is the same as the unsmoothed classifier  $f$  and the certified radius for a point is the distance to the decision boundary [8]. To make the problem analytically tractable, we use the squared-loss at the lower-level i.e.,  $f(x) = wx + b$  and  $l(x, y; f) = (f(x) - y)^2$ . The bilevel problem for poisoning is as follows

$$\begin{aligned} \min_u \quad & \mathbb{E}_{P^-} [\max(\text{sign}(w^*)(-b^*/w^* - x), 0)] \\ \text{s.t.} \quad & -\epsilon \leq u_i - x_i^- \leq \epsilon, \quad \text{for } i = 1, \dots, n \\ & w^*, b^* = \arg \min_{w, b} \frac{1}{2n} \left[ \sum_{i=1}^n l(x_i^+, 1) + \sum_{i=1}^n l(u_i, -1) \right]. \end{aligned} \quad (3)$$

**Theorem 1.** *If the perturbation is large enough, i.e.,  $\epsilon \geq \frac{\sum_i x_i^+ - \sum_i x_i^-}{n}$  then there are two locally optimal solutions to (3) which are  $u_i = x_i^- - \epsilon$  (Case 1) and  $u_i = x_i^- + \epsilon$  (Case 2) for  $i = 1, \dots, n$ . Otherwise, there is a unique globally optimal solution  $u_i = x_i^- - \epsilon$  (Case 1) for  $i = 1, \dots, n$ .*

Thus, optimal poisoning is achieved by shifting points of the  $P^-$  class either towards left or right by the maximum amount  $\epsilon$  (Fig. 2 and Appendix D.2). Moreover, the effect of poisoning an  $\alpha$  fraction of points from the  $P^-$  class with maximum permissible perturbation  $\tilde{\epsilon}$  is same as that of poisoning all points of  $P^-$  class with  $\epsilon = \alpha \tilde{\epsilon}$  (Corollary 2 in Appendix A). Although a direct analysis is intractable for non-linear cases, we empirically observed that our attack moved the decision boundary of neural networks closer to the points of the target class as measured by the mean distance of points to the decision boundary of the smoothed classifier (Sec. 4.6).

## 4. Experiments

In this section we present the results of our PACD<sup>1</sup> attack on poisoning deep neural networks trained using methods that make the model certifiably robust to test-time attacks. All the results presented here are averaged over models trained with five random initialization. We report the average certified radius (ACR) as the average of the certified radius obtained from the RS based certification procedure of [8] for correctly classified points. Certified radius is zero for misclassified and abstained points. The approximate certified accuracy (ACA) is the fraction of points correctly classified by the smoothed classifier ( $\ell_2$  radius greater than

<sup>1</sup>The code is available at [https://github.com/akshaymehra24/poisoning\\_certified\\_defenses](https://github.com/akshaymehra24/poisoning_certified_defenses)

Table 2. Decrease in certified radius and certified accuracy of models trained with Gaussian augmentation [8] on poison data compared to those of models trained on clean and watermarked data.

	$\sigma$	Data	Certified robustness of target class	
			ACR	ACA(%)
MNIST	0.25	Clean	0.896±0.01	98.92±0.32
		Watermarked	0.908±0.01	99.24±0.29
		Poisoned	<b>0.325±0.10</b>	<b>71.96±8.28</b>
	0.5	Clean	1.481±0.02	99.16±0.34
		Watermarked	1.514±0.06	99.12±0.47
		Poisoned	<b>0.733±0.10</b>	<b>90.68±3.37</b>
	0.75	Clean	1.549±0.11	98.48±0.35
		Watermarked	1.566±0.06	98.36±0.39
		Poisoned	<b>0.698±0.13</b>	<b>84.92±5.14</b>
CIFAR10	0.25	Clean	0.521±0.05	85.76±3.31
		Watermarked	0.470±0.01	83.22±1.41
		Poisoned	<b>0.059±0.02</b>	<b>26.84±6.04</b>
	0.5	Clean	0.634±0.04	75.04±1.65
		Watermarked	0.611±0.18	74.01±9.22
		Poisoned	<b>0.221±0.04</b>	<b>42.28±6.01</b>

zero). All results are reported over 500 randomly sampled images from the target classes. We use the same value of  $\sigma$  for smoothing during attack, retraining and evaluation. We compare our results to watermarking [31] which has been used previously for clean label attacks (opacity 0.1 followed by clipping to make  $\ell_\infty$  distortion equal to  $\epsilon$ ), and show that poison data generated using the bilevel optimization is significantly better at reducing the average certified radius.

We use our attack to poison MNIST and CIFAR10 dataset and use ApproxGrad to solve the bilevel optimization. The time complexity for ApproxGrad is  $O(VT)$  where  $V$  are the number of parameters in the machine learning model and  $T$  is the number of lower-level updates. For datasets like Imagenet where the optimization must be performed over a very large number of batches, obtaining the solution to bilevel problems becomes computationally hard. Due to this bottleneck we leave the problem of poisoning Imagenet for future work. For the experiments with MNIST we randomly selected digit 8 and for CIFAR10 the class ‘‘Ship’’ as the target class for the attacker. The attack results for other target classes are similar and are presented in the Appendix D. To ensure that the attack points satisfy the clean label constraint, the maximum permissible  $\ell_\infty$  distortion is bounded by  $\epsilon = 0.1$  for MNIST and  $\epsilon = 0.03$  for CIFAR10 which is similar to the value used to generate imperceptible adversarial examples in previous works [23, 12]. We used convolutional neural networks for our experiments on MNIST and Resnet-20 model for our experiments with CIFAR10. Model architectures, hyperparameters, generated attack examples (Fig. 5 in Appendix), and additional results on transferability of our poisoned samples to models with different architectures are presented in Appendix D. Since

Table 3. Decrease in certified radius and certified accuracy of models trained with MACER [36] on poison data compared to those of models trained on clean and watermarked data.

	$\sigma$	Data	Certified robustness of target class	
			ACR	ACA(%)
MNIST	0.25	Clean	0.915±0.01	99.64±0.21
		Watermarked	0.894±0.01	98.84±0.53
		Poisoned	<b>0.431±0.13</b>	<b>79.81±9.26</b>
	0.5	Clean	1.484±0.11	98.56±0.41
		Watermarked	1.475±0.08	98.68±0.39
		Poisoned	<b>0.685±0.16</b>	<b>84.36±6.17</b>
	0.75	Clean	1.353±0.13	93.81±2.08
		Watermarked	1.415±0.11	94.52±1.58
		Poisoned	<b>1.008±0.19</b>	<b>88.41±4.64</b>
CIFAR10	0.25	Clean	0.593±0.05	83.84±2.26
		Watermarked	0.486±0.04	77.01±0.21
		Poisoned	<b>0.379±0.11</b>	<b>72.41±9.79</b>
	0.5	Clean	0.759±0.11	72.92±5.06
		Watermarked	0.811±0.10	75.66±2.99
		Poisoned	<b>0.521±0.11</b>	<b>65.24±6.55</b>

models trained with standard training do not achieve high certified radius [8], we considered poisoning models trained with methods that improve the certified robustness guarantees to test time attacks. For comparison, ACR on the target class ‘‘Ship’’ with Resnet-20 trained with standard training on clean CIFAR10 dataset is close to zero whereas for the same model trained with GA ( $\sigma = 0.25$ ) ACR is close to 0.5. Finally, we show that our attack can withstand the use of weight regularization, which has been shown to be effective at mitigating the effect of poisoning attacks [6]. The results of this experiment are present in Appendix D.5.

#### 4.1. Poisoning Gaussian data augmentation [8]

Here we show the effectiveness of our attack at compromise the certified robustness guarantees obtained with RS on a model trained using the GA. The results of the attack, present in Table 2, show a significant decrease in the ACR and the certified accuracy of the target class of the model trained on poisoned data compared to the model trained on clean and watermarked data. Since the certified radius and certified accuracy are correlated, our poisoning attack which targets the reduction of certified radius (upper-level problem in Eq. (2)) also causes a decrease in the certified accuracy. Significant degradation in average certified radius from 0.52 to 0.06 on CIFAR10 with imperceptibly distorted poison data shows the extreme vulnerability of GA to poisoning.

#### 4.2. Poisoning MACER [36]

Here we use the bilevel formulation in Eq. (2) with  $\mathcal{L}_{\text{macer}}$  loss in the lower-level and generate poison data to reduce the certification guarantees of models trained with MACER. The poison data is generated with  $k = 2$ , where  $k$  are the number of noisy images used per training point to ease the

Table 4. Decrease in certified radius and certified accuracy of models trained with SmoothAdv [29] on poison data compared to those of models trained on clean and watermarked data.

	$\sigma$	Data	Certified robustness of target class	
			ACR	ACA(%)
MNIST	0.25	Clean	0.896±0.01	99.16±0.45
		Watermarked	0.906±0.01	99.28±0.16
		Poisoned	<b>0.672±0.04</b>	<b>93.21±1.92</b>
	0.5	Clean	1.408±0.05	99.21±0.25
		Watermarked	1.401±0.02	98.01±0.18
		Poisoned	<b>1.037±0.06</b>	<b>93.81±1.31</b>
0.75	Clean	1.262±0.05	95.68±0.47	
	Watermarked	1.433±0.03	97.21±0.13	
	Poisoned	<b>0.924±0.06</b>	<b>88.88±0.18</b>	
CIFAR10	0.25	Clean	0.504±0.02	78.76±0.81
		Watermarked	0.441±0.02	70.16±2.12
		Poisoned	<b>0.271±0.02</b>	<b>55.78±0.96</b>
	0.5	Clean	0.479±0.07	65.84±4.81
		Watermarked	0.473±0.02	62.51±2.12
		Poisoned	<b>0.277±0.02</b>	<b>49.11±3.19</b>

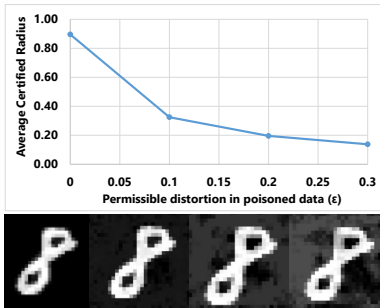


Figure 3. (Upper) ACR degrades more if larger perturbations are permitted to create poison data. But larger perturbation makes the poison points visibly distorted making them easier to detect with inspection (Lower). Poison data are generated with  $\epsilon \in \{0, 0.1, 0.2, 0.3\}$ . We have used  $\epsilon = 0.1$  for our attacks.

bilevel optimization. However, during retraining  $k = 16$  is used, which is similar to the one used in the original work [36]. The ACR obtained using MACER is higher than that achievable using Gaussian augmentation based training consistent with [36]. However, our data poisoning attack is still able to reduce the average certified radius of the method by more than 30% (Table 3) even though the attack is evaluated against a much stronger defense ( $k = 16$  for retraining compared to  $k = 2$  for poisoning) than what the poison data was optimized against. This shows that the use of a larger number of noisy samples ( $k$ ) cannot eliminate the effect of the attack, emphasising the importance of the threat posed by data poisoning.

### 4.3. Poisoning SmoothAdv [29]

Here we present the results of our attack on models trained with SmoothAdv. To yield model with high cer-

Table 5. Decrease in the mean  $\ell_2$ -distortion needed to generate adversarial examples using PGD attack against the smooth classifier. This shows that the decision boundary of the smooth classifier is closer to the clean test points of the target class after poisoning.

	$\sigma$	Clean data	Poisoned data
MNIST	0.25	3.271±0.10	<b>1.339±0.16</b>
	0.5	3.637±0.15	<b>2.170±0.09</b>
	0.75	3.961±0.18	<b>2.213±0.31</b>
CIFAR10	0.25	1.754±0.17	<b>0.132±0.04</b>
	0.5	1.996±0.09	<b>0.367±0.06</b>

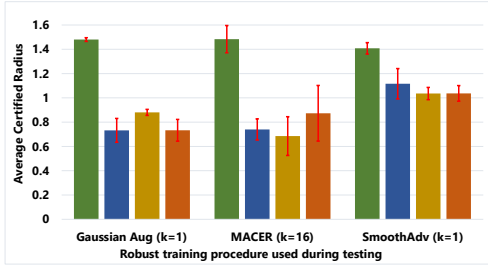
tified robustness this training method requires the model parameters to be optimized using adversarial training of the smoothed classifier. We used 2 step PGD attack to obtain adversarial example of each point in the batch. We used a single noisy instance of the adversarial example while doing adversarial training. Although, using larger  $k$  makes the certified robustness guarantees better, we used  $k = 1$  to save the computational time required for adversarial training. For bilevel training, we followed the similar procedure to generate adversarial examples for clean and poison data. The adversarial examples are then used as the data for the lower-level problem of Eq. (2) to do GA training for optimizing the network parameters. The batch of adversarial examples are recomputed against the updated model after each step of bilevel training. Note that this is an approximation to the actual solution of the minimax problem that has to be solved in the lower-level for generating poison data against SmoothAdv. However, the effectiveness of the attack (results in Table 4) suggests that our approximation works well in practice and certified robustness guarantees achieved from SmoothAdv can be degraded by poisoning.

### 4.4. Effect of the imperceptibility constraint

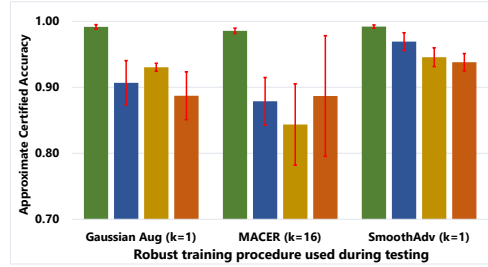
Here we evaluate the effect of using different values of the perturbation strength  $\epsilon$  which controls the maximum permissible distortion in Eq. (2). We use  $\sigma = 0.25$  for smoothing and GA based training to generate and evaluate the attack. The results are summarized in Fig. 3, which show that the ACR of the target class decreases as  $\epsilon$  increases rendering certification guarantees useless. This is expected since larger  $\epsilon$  creates a larger distribution shift among the target class data in the training and the test sets. However, larger permissible distortion to the data make the attack easier to detect by inspection. This is not desirable from an attacker’s perspective who wants to evade detection.

### 4.5. Transferability of poisoned data

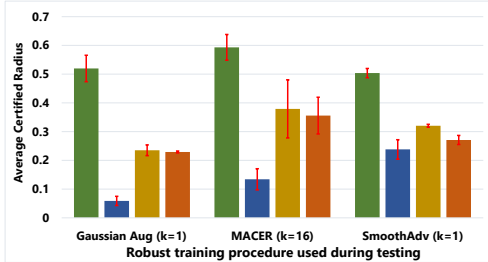
Here we report the performance of the models trained on the poison data using different training procedures than the one assumed by the attacker for crafting poison data. We used  $k = 1$  and 2 steps of PGD attack to generate adversarial examples for SmoothAdv and  $k = 16$  for MACER during



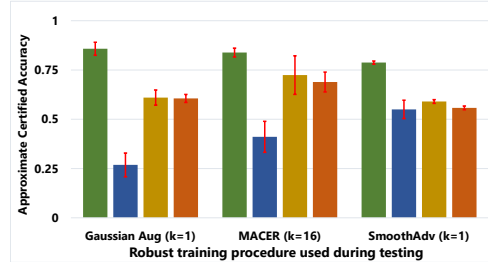
(a) Average certified radius of digit 8 in MNIST



(b) Approximate certified accuracy of digit 8 in MNIST



(c) Average certified radius of "Ship" in CIFAR10



(d) Approximate certified accuracy of "Ship" in CIFAR10

■ Clean data ■ Poison optimized against Gaussian Aug (k=1) ■ Poison optimized against MACER (k=2) ■ Poison optimized against SmoothAdv (k=1)

Figure 4. Successful transferability of our poisoned data when victim uses a training procedure different than the one used by the attacker to optimize the poison data.  $\sigma = 0.5$  and  $\sigma = 0.25$  are used for smoothing during training and evaluation for MNIST and CIFAR10.

retraining. The poison data generated against MACER was optimized using  $k = 2$ . The results are summarized in Fig. 4, which show that poisoned data optimized against any robust training procedure causes significant reduction in the certified robustness of models trained with a different training methods. Interestingly, poisoned data optimized against GA is extremely effective against other methods, considering the fact that it is the simplest of the three methods. The successful transferability of the poisoned data across different training methods shows how brittle these methods can be when faced with a poisoned dataset. We observe a similar success in transferability of the poison data to models with different architectures (Appendix D.4).

#### 4.6. Empirical robustness of poisoned classifiers

Finally, we report the empirical robustness of the smoothed classifier where the base classifier is trained on clean and poisoned data using GA. The poisoned data is generated against GA training in the lower-level as in Eq. (2). We report the mean  $\ell_2$ -distortion required to generate an adversarial example using the PGD attack [29] against the smoothed classifier using 200 and 100 randomly sampled test points of the target class from MNIST and CIFAR10, respectively, in Table 5. We observe that our poisoning leads to a decrease in the empirical robustness of the smoothed classifier on clean test data. This backs up our hypothesis that the decision boundary of the smooth classifier must be changed to reduce the certified radius in nonlinear classifiers, similar to linear classifiers (Fig. 2).

## 5. Conclusion

Certified robustness has emerged as a gold standard to gauge with certainty the susceptibility of machine learning models to test-time attacks. In this work, we showed that these guarantees can be rendered ineffective by our bilevel optimization based data poisoning attack that adds imperceptible perturbations to the points of the target class. Unlike previous data poisoning attacks, our attack can reduce the ACR of an entire target class and is even effective against models trained using training methods that have been shown to improve certified robustness. Our results suggests that data quality is a crucial factor in achieving high certified robustness guarantees but is overlooked by current approaches.

## 6. Acknowledgments

This work was supported by the NSF EPSCoR-Louisiana Materials Design Alliance (LAMDA) program #OIA-1946231 and by LLNL Laboratory Directed Research and Development project 20-ER-014 (LLNL-CONF-817233). This work was performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract No. DE-AC52-07NA27344, Lawrence Livermore National Security, LLC.<sup>2</sup>

<sup>2</sup>The views and opinions of the authors do not necessarily reflect those of the U.S. government or Lawrence Livermore National Security, LLC neither of whom nor any of their employees make any endorsements, express or implied warranties or representations or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of the information contained herein.



## References

- [1] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:1802.00420*, 2018. [1](#)
- [2] Jonathan F Bard. *Practical bilevel optimization: algorithms and applications*, volume 30. Springer Science & Business Media, 2013. [3](#)
- [3] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. *arXiv preprint arXiv:1206.6389*, 2012. [1](#)
- [4] Saikiran Bulusu, Bhavya Kailkhura, Bo Li, Pramod K Varshney, and Dawn Song. Anomalous example detection in deep learning: A survey. *IEEE Access*, 8:132330–132347, 2020. [1](#)
- [5] Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pages 3–14, 2017. [1](#)
- [6] Javier Carnerero-Cano, Luis Muñoz-González, Philippa Spencer, and Emil C Lupu. Regularisation can mitigate poisoning attacks: A novel analysis based on multiobjective bilevel optimisation. *arXiv preprint arXiv:2003.00040*, 2020. [3](#), [6](#), [14](#)
- [7] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017. [1](#)
- [8] Jeremy M Cohen, Elan Rosenfeld, and J Zico Kolter. Certified adversarial robustness via randomized smoothing. *arXiv preprint arXiv:1902.02918*, 2019. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [16](#)
- [9] Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pages 318–326, 2012. [3](#), [12](#)
- [10] Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and reverse gradient-based hyperparameter optimization. In *International Conference on Machine Learning*, pages 1165–1173, 2017. [3](#), [12](#)
- [11] Amin Ghiasi, Ali Shafahi, and Tom Goldstein. Breaking certified defenses: Semantic adversarial examples with spoofed robustness certificates. *arXiv preprint arXiv:2003.08937*, 2020. [3](#)
- [12] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. [6](#)
- [13] Sven Gowal, Krishnamurthy Dvijotham, Robert Stanforth, Rudy Bunel, Chongli Qin, Jonathan Uesato, Relja Arandjelovic, Timothy Mann, and Pushmeet Kohli. On the effectiveness of interval bound propagation for training verifiably robust models. *arXiv preprint arXiv:1810.12715*, 2018. [1](#)
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. [13](#)
- [15] Po-Sen Huang, Robert Stanforth, Johannes Welbl, Chris Dyer, Dani Yogatama, Sven Gowal, Krishnamurthy Dvijotham, and Pushmeet Kohli. Achieving verified robustness to symbol substitutions via interval bound propagation. *arXiv preprint arXiv:1909.01492*, 2019. [1](#)
- [16] W Ronny Huang, Jonas Geiping, Liam Fowl, Gavin Taylor, and Tom Goldstein. Metapoisson: Practical general-purpose clean-label data poisoning. *arXiv preprint arXiv:2004.00225*, 2020. [1](#), [3](#), [5](#), [15](#)
- [17] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. In *2018 IEEE Symposium on Security and Privacy (SP)*, pages 19–35. IEEE, 2018. [1](#)
- [18] Yujie Ji, Xinyang Zhang, and Ting Wang. Backdoor attacks against learning systems. In *2017 IEEE Conference on Communications and Network Security (CNS)*, pages 1–9. IEEE, 2017. [1](#)
- [19] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. *arXiv preprint arXiv:1703.04730*, 2017. [1](#)
- [20] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 656–672. IEEE, 2019. [1](#)
- [21] Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive noise. In *Advances in Neural Information Processing Systems*, pages 9464–9474, 2019. [1](#)
- [22] Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *International Conference on Machine Learning*, pages 2113–2122, 2015. [3](#), [12](#)
- [23] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017. [6](#)
- [24] Akshay Mehra and Jihun Hamm. Penalty method for inversion-free deep bilevel optimization. *arXiv preprint arXiv:1911.03432*, 2019. [3](#), [12](#)
- [25] Shike Mei and Xiaojin Zhu. Using machine teaching to identify optimal training-set attacks on machine learners. [3](#)
- [26] Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wonggrassamee, Emil C Lupu, and Fabio Roli. Towards poisoning of deep learning algorithms with back-gradient optimization. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pages 27–38, 2017. [1](#), [3](#)
- [27] Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In *International conference on machine learning*, pages 737–746, 2016. [3](#), [12](#)
- [28] Aditi Raghunathan, Jacob Steinhardt, and Percy S Liang. Semidefinite relaxations for certifying robustness to adversarial examples. In *Advances in Neural Information Processing Systems*, pages 10877–10887, 2018. [1](#)
- [29] Hadi Salman, Jerry Li, Ilya Razenshteyn, Pengchuan Zhang, Huan Zhang, Sebastien Bubeck, and Greg Yang. Provably robust deep learning via adversarially trained smoothed classifiers. In *Advances in Neural Information Processing Systems*, pages 11292–11303, 2019. [1](#), [2](#), [3](#), [4](#), [5](#), [7](#), [8](#), [16](#)
- [30] Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation for bilevel optimization. *arXiv preprint arXiv:1810.10667*, 2018. [3](#), [12](#)
- [31] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Su-

- ciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *Advances in Neural Information Processing Systems*, pages 6103–6113, 2018. [1](#), [5](#), [6](#)
- [32] Jacob Steinhardt, Pang Wei W Koh, and Percy S Liang. Certified defenses for data poisoning attacks. In *Advances in neural information processing systems*, pages 3517–3529, 2017. [1](#)
- [33] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Clean-label backdoor attacks. 2018. [1](#)
- [34] Jonathan Uesato, Brendan O’Donoghue, Aaron van den Oord, and Pushmeet Kohli. Adversarial risk and the dangers of evaluating against weak attacks. *arXiv preprint arXiv:1802.05666*, 2018. [1](#)
- [35] Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya Kailkhura, Xue Lin, and Cho-Jui Hsieh. Automatic perturbation analysis for scalable certified robustness and beyond. *Advances in Neural Information Processing Systems*, 33, 2020. [1](#)
- [36] Runtian Zhai, Chen Dan, Di He, Huan Zhang, Boqing Gong, Pradeep Ravikumar, Cho-Jui Hsieh, and Liwei Wang. Macer: Attack-free and scalable robust training via maximizing certified radius. *arXiv preprint arXiv:2001.02378*, 2020. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)
- [37] Chen Zhu, W Ronny Huang, Ali Shafahi, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein. Transferable clean-label poisoning attacks on deep neural nets. *arXiv preprint arXiv:1905.05897*, 2019. [1](#)

# Appendix

We present the proof of Theorem 1 and Corollary 2 in App. A followed by a review of bilevel optimization in App. B and our attack algorithm in App. C. In App. D, we present the results of additional experiments on poisoning different classes in the dataset, successful transferability of our poisoning attack to deeper models, performance of the attack when targeting a single test point and effect of using weight regularization on the attack success. We conclude in App. E by providing details of the hyperparameters and models architectures used in the experiments.

## A. Proofs

**Theorem 1.** *If the perturbation is large enough, i.e.,  $\epsilon \geq \frac{\sum_i x_i^+ - \sum_i x_i^-}{n}$  then there are two locally optimal solutions to (3) which are  $u_i = x_i^- - \epsilon$  (Case 1) and  $u_i = x_i^- + \epsilon$  (Case 2) for  $i = 1, \dots, n$ . Otherwise, there is a unique globally optimal solution which is  $u_i = x_i^- - \epsilon$  (Case 1) for  $i = 1, \dots, n$ .*

*Proof.* Let  $t = -\frac{b}{w}$  be the threshold of the linear classifier. Also let  $\Phi(t) := \int_{-\infty}^t P_-(x) dx$  and  $\Psi(t) := \int_{-\infty}^t x P_-(x) dx$ . There are two cases to consider.

**Case 1** ( $w > 0$ ): The upper-level cost function is

$$f(t) = \int_{-\infty}^t (t-x)P_-(x) dx = t\Phi(t) - \Psi(t)$$

Note that the range  $[-\infty, t]$  is where classification is correct for the test data. (Certified radius is 0 for misclassified points by definition.)

The closed-form solution of the lower-level problem gives us  $t = -\frac{b}{w} = \frac{\sum_i u_i + \sum_i x_i^+}{2n}$ , and therefore the perturbation bound  $|u_i - x_i^-| \leq \epsilon$  implies  $\sum_i x_i^- - n\epsilon \leq \sum_i u_i \leq \sum_i x_i^- + n\epsilon$  and therefore

$$-\frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n} \leq t \leq \frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n}.$$

Also, the assumption  $w > 0$  poses another constraint:

$w \propto \sum_i x_i^+ - \sum_i u_i > 0$  and therefore

$t = \frac{\sum_i u_i + \sum_i x_i^+}{2n} \leq \frac{\sum_i x_i^+}{n}$ . The upper-level problem is therefore

$$\begin{aligned} \min_t f(t) &= t\Phi(t) - \Psi(t) \quad \text{s.t.} \\ -\frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n} &\leq t \leq \frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n} \\ \text{and } t &\leq \frac{\sum_i x_i^+}{n}. \end{aligned}$$

Since  $f$  is non-decreasing (i.e.,  $f'(t) = \Phi(t) + tP_-(t) - tP_-(t) \geq 0$ ), the minimum is achieved at the left-most

boundary  $t = -\frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n}$  which corresponds to  $u_i = x_i^- - \epsilon$ ,  $i = 1, \dots, n$ .

**Case 2** ( $w < 0$ ): The upper-level cost function is now

$$f(t) = \int_t^{\infty} (-t+x)P_-(x) dx = -t(1-\Phi(t)) + (1-\Psi(t)),$$

which is non-increasing (i.e.,  $f'(t) = -(1-\Phi) + tP_- - tP_- \leq 0$ ) and has the constraints:

$$-\frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n} \leq t \leq \frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n}.$$

and

$$t = \frac{\sum_i u_i + \sum_i x_i^+}{2n} \geq \frac{\sum_i x_i^+}{n}.$$

For the solution to be feasible, it is required that  $\frac{\sum_i x_i^+}{n} \leq \frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n}$ , that is  $\frac{\epsilon}{2} \geq \frac{\sum_i x_i^+ - \sum_i x_i^-}{2n}$  (remember the assumption  $\frac{\sum_i x_i^-}{n} \leq \frac{\sum_i x_i^+}{n}$ ). Therefore if the perturbation is large enough, i.e.,  $\epsilon \geq \frac{\sum_i x_i^+ - \sum_i x_i^-}{n}$  holds, then the minimum is achieved at the right-most boundary  $t = \frac{\epsilon}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n}$  which corresponds to  $u_i = x_i^- + \epsilon$ ,  $i = 1, \dots, n$ .  $\square$

**Corollary 2.** *If the perturbation is large enough, i.e.,  $\epsilon \geq \frac{\sum_i x_i^+ - \sum_i x_i^-}{n}$  then the reduction in the ACR of the target class, by poisoning an  $\alpha$  portion ( $\alpha \in [0, 1]$ ) of the target class, with maximum perturbation  $\tilde{\epsilon}$  using Eq. (3) is same as that achieved by poisoning the entire target class with  $\epsilon = \alpha\tilde{\epsilon}$ .*

*Proof.* To obtain the effect of poisoning an  $\alpha$  portion ( $\alpha \in [0, 1]$ ) of the target class, with maximum perturbation  $\tilde{\epsilon}$ , we follow the proof of Theorem 1, with the change that the solution to the lower-level problem is  $t = -\frac{b}{w} = \frac{\sum_{i=0}^{\alpha n} u_i + \sum_{i=0}^{(1-\alpha)n} x_i^- + \sum_i x_i^+}{2n}$ .

Thus the solutions to the upper-level problem are

$t = \frac{-\alpha\tilde{\epsilon}}{2} + \frac{\sum_i x_i^+ + \sum_{i=0}^{\alpha n} x_i^- + \sum_{i=0}^{(1-\alpha)n} x_i^-}{2n}$  which corresponds to  $u_i = x_i^- - \tilde{\epsilon}$ ,  $i = 1, \dots, \alpha n$  for Case 1 and

$t = \frac{\alpha\tilde{\epsilon}}{2} + \frac{\sum_i x_i^+ + \sum_{i=0}^{\alpha n} x_i^- + \sum_{i=0}^{(1-\alpha)n} x_i^-}{2n}$  which corresponds to  $u_i = x_i^- + \tilde{\epsilon}$ ,  $i = 1, \dots, \alpha n$  for Case 2.

The decision boundaries

$t = \frac{-\alpha\tilde{\epsilon}}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n}$  for Case 1 and

$t = \frac{\alpha\tilde{\epsilon}}{2} + \frac{\sum_i x_i^+ + \sum_i x_i^-}{2n}$  for Case 2 are the same boundaries as obtained by poisoning all points from the target class ( $x_i^-$ ) with  $\epsilon = \alpha\tilde{\epsilon}$   $\square$

## B. Review of bilevel optimization

A bilevel optimization problem is of the form  $\min_{u \in \mathcal{U}} \xi(u, v^*)$  s.t.  $v^* = \arg \min_{v \in \mathcal{V}(u)} \zeta(u, v)$ , where the upper-level problem is a minimization problem with  $v$  constrained to be the optimal solution to the lower-level problem. General bilevel problems are difficult to solve but if the solution to the lower-level problem can be computed in closed form then we can replace the lower-level problem with its solution, reducing the bilevel problem into a single level problem. We can then use the gradient-based methods to solve the single level problem. The total derivative  $\frac{d\xi}{du}(u, v^*(u))$  (hypergradient) using the chain rule is

$$\frac{d\xi}{du} = \nabla_u \xi + \frac{dv}{du} \cdot \nabla_v \xi.$$

Since  $\nabla_v \zeta = 0$  at  $v = v^*(u)$  and assuming  $\nabla_{vv}^2 \zeta$  is invertible we can compute  $\frac{dv}{du}$  using the implicit function theorem (this can be done even if the solution to lower-level problem can't be found in closed form) which gives

$$\frac{dv}{du} = -\nabla_{uv}^2 \zeta (\nabla_{vv}^2 \zeta)^{-1}.$$

Thus the hypergradient is

$$\frac{d\xi}{du} = \nabla_u \xi - \nabla_{uv}^2 \zeta (\nabla_{vv}^2 \zeta)^{-1} \nabla_v \xi \text{ at } (u, v^*(u)).$$

Since computation of  $(\nabla_{vv}^2 \zeta)^{-1}$  is difficult, [9, 27] proposed to instead approximate the solution to  $q = (\nabla_{vv}^2 \zeta)^{-1} \nabla_v \xi$  by approximately solving the linear system of equations  $\nabla_{vv}^2 \zeta \cdot q \approx \nabla_v \xi$ . This can be done by minimizing  $\|\nabla_{vv}^2 \zeta \cdot q - \nabla_v \xi\|$  using any iterative solver. Other methods for solving the bilevel optimization problems include using forward/reverse mode differentiation [10, 22, 30] to approximate the inverse and penalty method [24] to solve the single level problem as a constrained minimization problem.

## C. Attack algorithm

Alg. 1 shows the complete algorithm used to generate the poisoning attack when RS is used for certification and models are trained using GA. The algorithm relies on ApproxGrad (Alg. 2) to solve the bilevel optimization problem. The upper-level cost is a differentiable function that approximates the certified radius of the hard smooth classifier using a soft smooth classifier. The hyperparameter  $\alpha$  is the inverse temperature parameter of softmax. As  $\alpha \rightarrow \infty$ , softmax converges to argmax almost everywhere. As a result  $\tilde{g}_\theta$  converges to  $g_\theta$  almost everywhere and thus soft randomized smoothing converges to hard randomized smoothing almost everywhere. Although, in this work we considered RS as the procedure for certification (due to its scalability to large models and datasets), any other certification procedure can

---

### Algorithm 2 Algorithm for ApproxGrad

---

Input:  $\xi, \zeta, M, T_1, T_2, \epsilon, u_{base}, \{\tau_m = 0.1\}, \{\rho_{m,t_1} = 0.001\}, \{\beta_{m,t_2} = 0.001\}$

Output:  $(u_K)$

Initialize  $u_0, v_0$  randomly

Begin

**for**  $m = 0, \dots, M-1$  **do**

    # Approximately solve the lower-level problem

**for**  $t = 0, \dots, T_1-1$  **do**

$v_{t+1} \leftarrow v_t - \rho_{m,t_1} \nabla_v \zeta$

**end for**

    # Approximately solve the linear system

    #  $\nabla_{vv}^2 \zeta \cdot q_k = \nabla_v \xi$

**for**  $t = 0, \dots, T_2-1$  **do**

$q_{t+1} \leftarrow q_t - \beta_{m,t_2} \nabla_q (\|\nabla_{vv}^2 \zeta \cdot q_m - \nabla_v \xi\|)$

**end for**

    # Compute the approximate Hypergradient

$p_m = \nabla_u \xi - \nabla_{uv}^2 \zeta \cdot q_{T_2}$

    # Update  $u_m$  and use projection for the

    # upper-level constraint

$u_{m+1} = P(u_m - \tau_m p_m, \epsilon, u_{base})$

**end for**

---

be used as the upper-level cost as long as its differentiable. Moreover, Alg. 1 uses  $\mathcal{L}_{\text{GaussAug}}$  in the lower-level to train the model, but like the case with upper-level cost any other loss function can be used to obtain the model parameters. This flexibility of our method allows us to generate poison data against MACER and SmoothAdv using their loss functions in the lower-level.

### C.1. ApproxGrad

For an unconstrained bilevel problem of the form  $\min_u \xi(u, v^*)$  s.t.  $v^* = \arg \min_v \zeta(u, v)$ , if  $\zeta(u, v)$  is strongly convex then we can replace the lower-level problem with its necessary condition for optimality and write the bilevel problem as the following single level problem  $\min_u \xi(u, v^*)$  s.t.  $\nabla_v \zeta(u, v) = 0$ . Assuming  $\nabla_{vv}^2 \zeta$  is invertible everywhere we can compute the hypergradient at the point  $(u, v^*(u))$  as  $\frac{d\xi}{du} = \nabla_u \xi - \nabla_{uv}^2 \zeta (\nabla_{vv}^2 \zeta)^{-1} \nabla_v \xi$ .

The ApproxGrad algorithm approximates the Hessian-inverse vector product by approximately solving a system of linear equation using an iterative solver such as gradient descent or conjugate gradient method. In this work we use Adam optimizer to solve this system. Since our problem for data poisoning in Eq. (2) involves a constraint in the upper-level we use projection to enforce the constraint. The

full algorithm for solving the bilevel optimization problem using ApproxGrad is present in Alg. 2. For our attack the lower-level problem involves a deep neural network, which can have multiple local minima and thus optimizing against a single local minima in the bilevel problem is not ideal. To overcome this problem we reinitialize the lower-level variable  $v$  after few upper-level iterations to prevent the poisoning points from overfitting to a particular local minima. Empirically, this helps us find poisoning points that remain effective even after the model is retrained from scratch making them generalize to different initialization of the neural network.

## D. Additional experiments

### D.1. Comparison with standard data poisoning

The standard data poisoning attack creates poison data so that the accuracy of the victim’s model trained on it is significantly lower than the accuracy attainable with training on clean data. The bilevel optimization problem for this attack is as follows.

$$\begin{aligned} & \min_u \mathcal{L}_{\text{standard}}(\mathcal{D}^{\text{val}}) \\ & \text{s.t. } \|\delta_i\|_{\infty} \leq \epsilon, \quad i = 1, \dots, n, \quad \text{and} \quad (4) \\ \theta^* & = \arg \min_{\theta} \mathcal{L}_{\text{standard}}(\mathcal{D}^{\text{clean}} \cup \mathcal{D}^{\text{poison}}; \theta). \end{aligned}$$

Here  $\mathcal{L}_{\text{standard}}(\mathcal{D}; \theta) = \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} l_{ce}(x_i, y_i; \theta)$ , where  $l_{ce}$  is the cross entropy loss. We used this formulation to generate the poisoned dataset for reporting the results in Table 1 with  $\epsilon = 0.1$  for MNIST and  $\epsilon = 0.03$  for CIFAR10. The attack modifies all the points in the target classes. Specifically, our attack targets misclassification of the digit 8 in MNIST and class “Ship” in CIFAR10. The poisoned dataset obtained after solving the bilevel optimization was then used to train five models starting from random initializations with different training procedures. The results of which are reported in Table 1. As expected the models trained with standard training on the poisoned data perform the worst in terms of accuracy since the attack was optimized against standard training. However, the generated poison data has little to no effect when a training procedures that improves certified adversarial robustness is used. This shows that the effect of standard data poisoning can easily be nullified if a victim trains the model with a these training procedures. This gives a false sense of security of the models trained with certified defenses to data poisoning attacks. Thus, in this work we study the effect of poisoning on training procedures meant to improve certified adversarial robustness and show that their guarantees become meaningless when the dataset is poisoned.

Table 6. Degradation of certified adversarial robustness of logistic regression trained with GA on a toy 2-isotropic Gaussians dataset.

$\sigma$	Certified Robustness on clean data		Certified Robustness on poisoned data	
	ACR	ACA(%)	ACR	ACA(%)
0.25	0.4047	90.00	0.3585	88.00
0.50	0.4139	90.00	0.3587	87.60
0.75	0.4123	90.00	0.3544	87.60

### D.2. Isotropic Gaussians

Here we validate the solution found by solving the bilevel optimization against the analytical solution of a toy problem. Consider a two-dimensional dataset comprising of points drawn from two isotropic Gaussian distributions. Let  $\mathbb{P}(x|y = -1) = \mathcal{N}(\mu_1, \sigma^2 I)$  and  $\mathbb{P}(x|y = 1) = \mathcal{N}(\mu_2, \sigma^2 I)$  and equal prior  $\mathbb{P}(y = 1) = \mathbb{P}(y = -1)$ . For a point  $x$ , the Bayes optimal classifier predicts  $y = 1$  if  $\mathbb{P}(y = 1|x) \geq \mathbb{P}(y = -1|x)$  and predicts  $y = -1$  otherwise. The decision boundary of the Bayes optimal classifier is given by  $(\mathbf{x} - \mu_1)^T(\mathbf{x} - \mu_1) = (\mathbf{x} - \mu_2)^T(\mathbf{x} - \mu_2)$ . This is also the decision boundary of the smoothed classifier. Assuming the attacker is poisoning the class with label  $-1$  and maximum permissible distortion is  $\epsilon$ , our analysis showed that maximum reduction in radius occurs if the entire distribution shifts by  $\epsilon$  i.e. the new mean of the class with label  $-1$  is  $\mu_1 - \epsilon$  and the decision boundary is  $(\mathbf{x} - (\mu_1 - \epsilon))^T(\mathbf{x} - (\mu_1 - \epsilon)) = (\mathbf{x} - \mu_2)^T(\mathbf{x} - \mu_2)$ . Since the test distribution is unchanged, the ACR for the test points with labels  $-1$  is reduced by  $\frac{\epsilon}{\sqrt{2}}$ . Using  $\mu_1 = 0.2, \mu_2 = 0.8, \sigma_1 = \sigma_2 = 0.3, \epsilon = 0.1$  and using logistic regression in the lower-level, analytically the certified radius must decrease from 0.4243 to 0.3546. The solution by solving the bilevel optimization numerically (Table 6) matches the analytic solution.

### D.3. Targeting other classes

In this section we present the results of our poisoning attack on different target classes where the models are trained using GA during poison generation and evaluation. Since MNIST and CIFAR10 both have 10 classes we create 10 poisoning sets each targeting a particular class. The results of retraining models from five random initializations on each of the 10 poisoning sets are summarized in Fig. 6. Reduction in average certified radius for all classes shows that an attacker can generate poison data to affect any class in the dataset.

### D.4. Transferability to different architectures

Here we present the results of transferability of the poisoned data generated against Resnet-20 targeting the class “Ship” to bigger models. In particular we present the results on Resnet-56 and Resnet-110 [14] models in Fig. 7. As seen from the results the poisoned data generated against



(a) Clean data (odd numbered rows) and poisoned data generated by our attack (even numbered rows) for all digits in MNIST



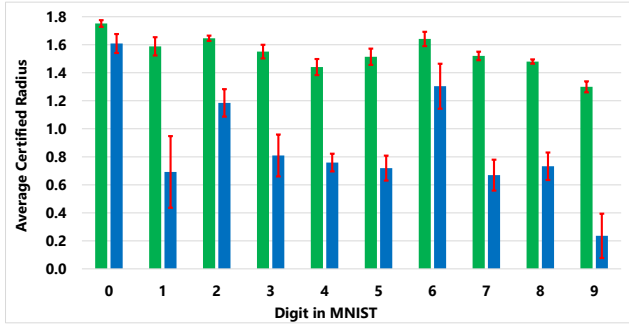
(b) Clean data (odd numbered rows) and poisoned data generated by our attack (even numbered rows) for all classes of CIFAR10

Figure 5. Imperceptibly distorted poison data generated by our algorithm against Gaussian augmented training which causes a significant reduction in the certified robustness guarantees of the models. The average certified radius and certified accuracy of models trained on clean and poisoned data are reported in App. D.3 and Fig. 6.

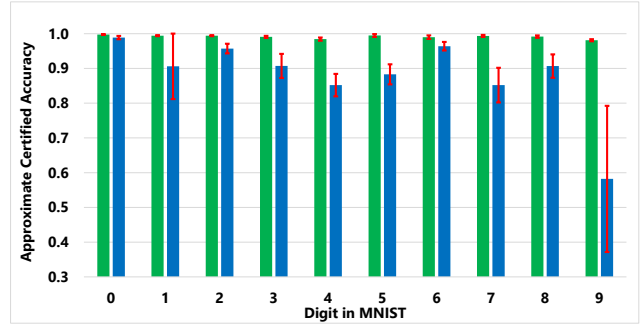
Resnet-20 is successful in reducing the certified radius of the target class even if the victim uses a larger model. We report the results of training the models on clean and poisoned data starting from three random initializations and certify using 500 randomly sampled points of the target class from the clean test set. Our results suggest that the poisoned data generated using our procedure are agnostic to the training procedure (Fig. 4), model (Fig. 7) and metric (RS or empirical robustness) used by the victim during evaluation highlighting the threat of data poisoning.

## D.5. Effect of weight regularization

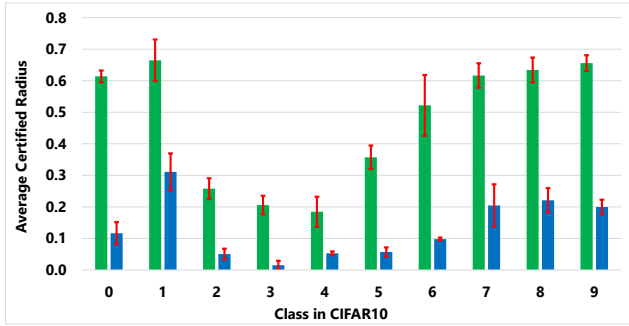
Previous works [6] have shown weight regularization to mitigate the effect of data poisoning attacks. Here, we evaluate the attack success when using different coefficients for weight regularization in standard and GA training. Results in Table 7 show that our attack significantly reduces the ACR of models, especially those trained without GA or weight regularization. Similar to previous works [6], we see that models trained with large regularization (without GA) are difficult to poison. This increased robustness, however comes at the cost



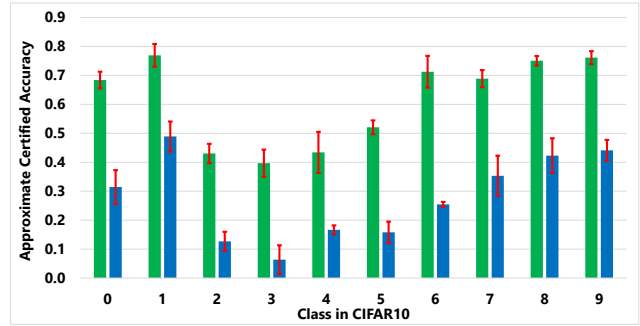
(a) Average certified radius of all digits in MNIST



(b) Approximate certified accuracy of all digits in MNIST



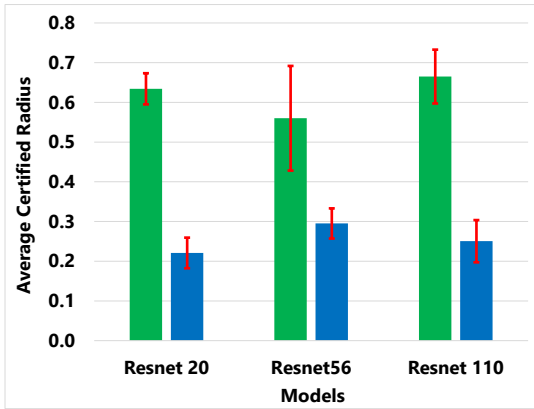
(c) Average certified radius of all classes in CIFAR10



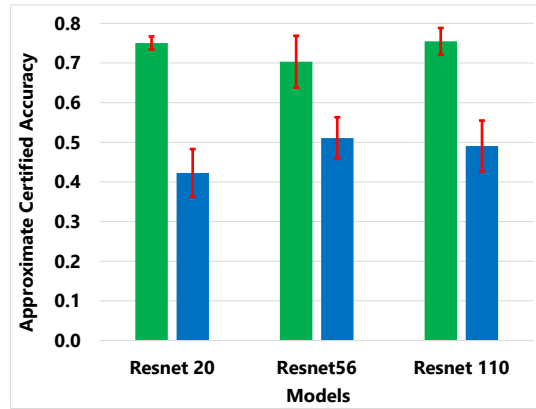
(d) Approximate certified accuracy of all classes in CIFAR10

■ Clean data ■ Poisoned data

Figure 6. Successful poisoning attack against all classes in MNIST and CIFAR10 dataset.



(a) Average certified radius of class "Ship" in CIFAR10



(b) Approximate certified accuracy of class "Ship" in CIFAR10

■ Clean data ■ Poisoned data

Figure 7. Successful transferability of poisoning attack against deeper models. The poison data is optimized against Resnet-20. GA with  $\sigma = 0.5$  is used for poison generation and evaluation.

of accuracy (target class accuracy of models trained with clean data drops from 99% to 92%), suggesting a trade-off. On the other hand, for models trained with GA, using large regularization leads to a significant drop in their certified robustness guarantees, even without poisoning (ACR for models trained on clean data drops from 1.48 to 0.85), there by making large regularization undesirable to use with GA.

This shows that our attack remains quite effective even when different amounts of weight regularization are used during model retraining.

### D.6. Attack success by poisoning 1% of the data

Similar to [16], where the effect of poisoning 1% of the training data is evaluated on a single test point, we randomly

Table 7. Effect of weight regularization and Gaussian data augmentation ( $\sigma=0.5$ ) on ACR of digit 8 of MNIST ( $\epsilon=0.1$ ). Mean and s.d. of 3 random initializations. Bold entries are reported in Table 2.

Training	Data	No Reg.	1E-4	1E-2	1E-1
Without GA	Clean	0.95±0.10	0.89±0.06	0.87±0.11	0.82±0.04
	Poisoned	0.01±0.01	0.03±0.05	0.37±0.06	0.68±0.05
With GA	Clean	<b>1.48±0.02</b>	1.49±0.03	1.29±0.15	0.85±0.09
	Poisoned	<b>0.73±0.10</b>	0.62±0.02	0.58±0.11	0.72±0.08

select 5 test images of the bird class and generate poison data to reduce their certified radius individually. Using Alg. 1, we poison 500 bird images (1% of CIFAR10) closest to the target, with  $\epsilon=0.06$ . Using 3 randomly initialized models trained with GA, our attack can reduce the certified radius of 5 targets from 0.63 to 0.26 on average.

## E. Details of the experiments

All codes are written in Python using Tensorflow/Keras, and were run on Intel Xeon(R) W-2123 CPU with 64 GB of RAM and dual NVIDIA TITAN RTX. Implementation and hyperparameters are described below.

### E.1. Data splits

For MNIST, we use 55000 points as the training data and 5000 points for validation data. We have roughly 500 points belonging to the target class in the validation set which is used in the upper-level problem of the bilevel optimization presented in Eq. (1). For CIFAR10, we use 45000 points as the training data and 5000 points for validation data. Similar to MNIST we have roughly 500 points belonging to the target class in the validation set. The test sets of both the datasets comprises of 10000 points. We use 500 randomly sampled points of the target class from the test set to report the results of certified and empirical robustness of the models trained on clean and poisoned data.

### E.2. Model Architecture

For the experiments on the MNIST dataset, our network consists of a convolution layer with kernel size of 5x5, 20 filters and ReLU activation, followed by a max pooling layer of size 2x2. This is followed by another convolution layer with 5x5 kernel, 50 filters and ReLU activation followed by similar max pooling and dropout layers. Then we have a fully connected layers with ReLU activation of size 500. Lastly, we have a softmax layer with 10 classes. The accuracy of the model on clean data when optimized with the Adam optimizer using a learning rate of 0.001 for 100 epochs with batch size of 200 is 99.3% (without GA). For the experiments on the CIFAR10 dataset, we use the Resnet-20 model. The accuracy of the model on clean data when optimized with the Adam optimizer using a learning rate of 0.001 for 100 epochs with batch size of 200 is 85% (without GA). For all CIFAR10 experiments except for the experiments with SmoothAdv, we

trained the models using data augmentation (random flipping and random cropping). We used the same parameters for training the models with different robust training procedures on clean and poisoned data.

### E.3. Hyperparameters

For experiments with MNIST we used  $\epsilon = 0.1, K = 20, \alpha = 16$ . The batch size used for lower-level training was 1000, of which 100 points belonged to the poisoned set (target class). The batch size for validation set was 100 which only consisted of points from the target class. The lower-level was trained using different training procedures on clean and poisoned data. For experiments with CIFAR10 we used  $\epsilon = 0.03, \lambda = 0.06, M = 20, \alpha = 16$ . The batch size used for lower-level training was 200, of which 20 points belonged to the poisoned set (target class). The batch size for validation set was 20 which only consisted of points from the target class. For training with GA the lower-level is trained with a single noisy image of the clean and poisoned dataset. The same setting is used while retraining. For generating poison data against MACER the lower-level is trained with  $K = 2, \lambda = 1, \gamma = 8$ . During retraining,  $K = 16, \lambda = 16, \gamma = 8$  are used for MNIST. For CIFAR10, we use  $K = 16, \gamma = 8$  and  $\lambda = 12$  for  $\sigma = 0.25$  and  $\lambda = 4$  for  $\sigma = 0.5$ . The hyperparameters during retraining are similar to the ones used in the original work. For Smoothadv,  $k = 1$  and 2-step PGD attack are used to generate adversarial examples of the smooth classifier. These adversarial examples along with GA are used to do adversarial training during poison generation and retraining.

In our experiments with generating poisoned data against GA and MACER we used  $P = 50, T_1 = T_2 = 10, \tau = 0.1, \rho = 0.001, \beta = 0.01$  for ApproxGrad. We used all the same hyperparameters for SmoothAdv except  $T_1 = 1$ . For certification we used the CERTIFY procedure of [8], with  $n_0 = 100, n = 100000, \alpha = 0.001$ . For measuring empirical robustness of the smoothed classifier, we used the mean  $\ell_2$  distortion required by PGD attack to generate an adversarial example as done in [29]. The attack is optimized for 100 iterations for different values of  $\ell_2$  distortion between (0.01, 10). We used 20 augmentations for each test point of MNIST and 10 for CIFAR10. To report the results for empirical robustness we record the minimum distortion for a successful attack for each test point.

For the watermarking baseline, we randomly selected an image (*other*) from the classes other than the target class and over-layed them on top of the target class images (*base*) with an opacity of  $\gamma = 0.1$  i.e. ( $poison\_image = \gamma \cdot other + (1 - \gamma) \cdot base$ ). We then clip the images to have  $\ell_\infty$  distortion of  $\epsilon$  to make our bilevel attack comparable in terms of maximum distortion.