

# Data-Dependent Randomized Smoothing

Motasem Alfarra<sup>1,\*</sup>

Adel Bibi<sup>2,\*</sup>

Philip H. S. Torr<sup>2</sup>

Bernard Ghanem<sup>1</sup>

<sup>1</sup> King Abdullah University of Science and Technology (KAUST), Saudi Arabia

<sup>2</sup>University of Oxford, United Kingdom

## Abstract

Randomized smoothing is a recent technique that achieves state-of-art performance in training certifiably robust deep neural networks. While the smoothing family of distributions is often connected to the choice of the norm used for certification, the parameters of these distributions are always set as global hyper parameters independent from the input data on which a network is certified. In this work, we revisit Gaussian randomized smoothing and show that the variance of the Gaussian distribution can be optimized at *each* input so as to maximize the certification radius for the construction of the smooth classifier. Since the data dependent classifier does not directly enjoy sound certification with existing approaches, we propose a memory-enhanced data dependent smooth classifier that is certifiable by construction. This new approach is generic, parameter-free, and easy to implement. In fact, we show that our data dependent framework can be seamlessly incorporated into 3 randomized smoothing approaches, leading to consistent improved certified accuracy. When this framework is used in the training routine of these approaches followed by a data dependent certification, we achieve 9% and 6% improvement over the certified accuracy of the strongest baseline for a radius of 0.5 on CIFAR10 and ImageNet.

## 1 INTRODUCTION

<sup>1</sup> Despite the success of Deep Neural Networks (DNNs) in various learning tasks [Krizhevsky et al., 2012, Long et al.,

\* Equal contribution.

Correspondence to: motasem.alfarra@kaust.edu.sa, adel.bibi@eng.ox.ac.uk

<sup>1</sup>Code: <https://github.com/MotasemAlfarra/DDRS>.

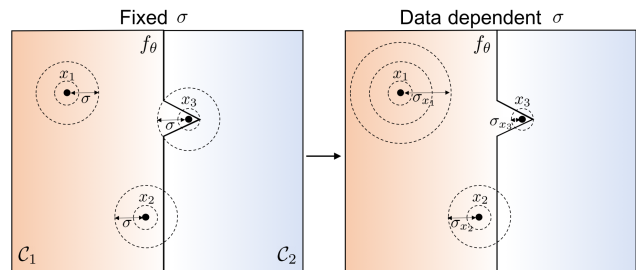


Figure 1: **From fixed to data dependent smoothing.** Using a fixed  $\sigma$  for all inputs to smooth  $f_\theta$  may under certify (results in smaller certification radius) inputs far from decision boundary *e.g.*  $x_1$ , decrease in prediction confidence as for  $x_2$  or produce incorrect predictions as for  $x_3$ . Thus, smoothing should vary per input (right figure) to alleviate the aforementioned issues.

2015], they were shown to be vulnerable to small carefully crafted adversarial perturbations [Goodfellow et al., 2015, Szegedy et al., 2013]. For a DNN  $f$  that correctly classifies an image  $x$ ,  $f$  can be fooled to produce an incorrect prediction for  $x + \eta$  even when the adversary  $\eta$  is so small that  $x$  and  $x + \eta$  are indistinguishable to the human eye. To circumvent this nuisance, there have been several works proposing heuristic training procedures to build networks that are *robust* against such perturbations [Cisse et al., 2017, Madry et al., 2018]. However, many of these works provided a false sense of security as they were subsequently broken, *i.e.* shown to be ineffective against stronger adversaries [Athalye et al., 2018, Tramer et al., 2020, Uesato et al., 2018]. This has inspired researchers to develop networks that are *certifiably robust*, *i.e.* networks that provably output constant predictions over a characterized region around every input. Among many certification methods, a probabilistic approach to certification called *randomized smoothing* has demonstrated impressive state-of-the-art certifiable robustness results [Cohen et al., 2019, Lecuyer et al., 2019, Li et al., 2019]. In a nutshell, given an input  $x$  and a base classifier  $f$ , *e.g.* a DNN, randomized smoothing constructs

a “smooth classifier”  $g(x) = \mathbb{E}_{\epsilon \sim \mathcal{D}} [f(x + \epsilon)]$  such that, and under some choices of  $\mathcal{D}$ ,  $g(x) = g(x + \delta) \forall \delta \in \mathcal{R}$ . As such,  $g$  is certifiable within the certification region  $\mathcal{R}$  characterized by  $x$  and the smoothing distribution  $\mathcal{D}$ . While there has been considerable progress in devising a notion of “optimal” smoothing distribution  $\mathcal{D}$  for when  $\mathcal{R}$  is characterized by an  $\ell_p$  certificate [Yang et al., 2020], a common trait among all works in the literature is that the choice of  $\mathcal{D}$  is independent from the input  $x$ . For example, one of the earliest works on randomized smoothing grants  $\ell_2$  certificates under  $\mathcal{D} = \mathcal{N}(0, \sigma^2 I)$ , where  $\sigma$  is a free parameter that is constant for all  $x$  [Cohen et al., 2019]. That is to say, the classifier  $f$  is smoothed to a classifier  $g$  uniformly (same variance  $\sigma^2$ ) over the entire input space of  $x$ . The choice of  $\sigma$  used for certification is often set either arbitrarily or via cross validation to obtain best certification results [Salman et al., 2019a]. We believe this is suboptimal and that  $\sigma$  should vary with the input  $x$  (data dependent), since using a fixed  $\sigma$  may under-certify inputs (*i.e.* the constructed smooth classifier  $g$  produces smaller certification radii), which are far from the decision boundaries as exemplified by  $x_1$  in Figure 1. Moreover, this fixed  $\sigma$  could be large for inputs  $x$  close to the decision boundaries resulting in a smooth classifier  $g$  that incorrectly classifies  $x$  (refer to  $x_3$  in Figure 1).

In this paper, we aim to introduce more structure to the smoothing distribution  $\mathcal{D}$  by rendering its parameters data dependent. That is to say, the base classifier  $f$  is smoothed with a family of smoothing distributions to produce:  $g(x) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma_x^2 I)} [f(x + \epsilon)]$ <sup>2</sup>. Note here that the variance of the Gaussian is now dependent on the data input  $x$ . Moreover, given that  $\sigma_x$  varies with  $x$ , classical randomized smoothing based certification does not apply directly. We propose a simple memory-based approach to certify the resultant data dependent smooth classifier  $g$ . We show that our memory-enhanced data dependent smooth classifier can boost certification performance of several randomized smoothing techniques. Our contributions can thus be summarized in three folds. (i) We propose a parameter free and generic framework that can easily turn several randomized smoothing techniques into their data dependent variants. In particular, given a network  $f$  and an input  $x$ , we propose to optimize the smoothing distribution parameters for every  $x$ , *e.g.*  $\sigma_x^*$ , so they maximize the certification radius. This choice of  $\sigma_x^*$  is then used to smooth  $f$  at  $x$  and construct a smoothed classifier  $g$ . Moreover, as the data dependent smooth classifier is not directly certifiable using Cohen et al. [2019] MCMC approaches, we propose a memory-enhanced data dependent smooth classifier for certification. (ii) We demonstrate the effectiveness of our memory-enhanced data dependent smoothing by showing that we can improve the certified accuracy of several models, specifically models trained with Gaussian augmentation

(COHEN) [Cohen et al., 2019], adversaries on the smoothed classifier (SMOOTHADV) [Salman et al., 2019a], and radius regularization (MACER) [Zhai et al., 2020] *without any model retraining*. We boost the certified accuracy of the best baseline by 5.4% on CIFAR10 and by 2.8% on ImageNet for  $\ell_2$  perturbations with less than 0.5 (=127/255) ball radius. (iii) We show that incorporating the proposed data dependent smoothing in the training pipeline of COHEN, SMOOTHADV and MACER can further boost results to get certified accuracies of 68.3% on CIFAR10 and 64.2% on ImageNet at  $\ell_2$  perturbations less than 0.25.

## 2 RELATED WORK

**Certified Defenses.** Certified defenses aim to guarantee that an adversary does not exist in a certain region around a given input. Certified defenses can be divided into exact [Cheng et al., 2017, Lomuscio and Maganti, 2017, Huang et al., 2017, Ehlers, 2017] and relaxed certification [Salman et al., 2019b, Wong and Kolter, 2018]. Generally, exact certification suffers from poor scalability with networks that are at most 3 hidden layers deep [Tjeng et al., 2019]. On the other hand, relaxed methods resolve this issue by aiming at finding an upper bound to the worst adversarial loss over all possible bounded perturbations around a given input [Weng et al., 2018]. However, the latter is too expensive for any mixed certification-training routine.

**Randomized Smoothing.** The earliest work on randomized smoothing [Lecuyer et al., 2019] was from a differential privacy perspective, where it was demonstrated that adding Laplacian noise enjoys an  $\ell_1$  certification radius in which the average classifier prediction under this noise is constant. This work was later followed by the tight  $\ell_2$  certificate radius for Gaussian smoothing [Cohen et al., 2019]. Since then, there has been a body of work on randomized smoothing with empirical defenses [Salman et al., 2019a] to certify black box classifiers [Salman et al., 2020]. Other works derived certification guarantees for  $\ell_1$  bounded [Teng et al., 2019],  $\ell_\infty$  bounded [Zhang et al., 2019], and  $\ell_0$  bounded [Levine and Feizi, 2020] perturbations. Even more recently, a novel framework that finds the optimal smoothing distribution for a given  $\ell_p$  norm [Yang et al., 2020] was proposed showing state-of-art certification results on  $\ell_1$  perturbations. We deviate from the common literature by introducing the notion of smoothing, particularly Gaussian smoothing for  $\ell_2$  perturbations, which varies depending on the input. In particular, since an input  $x$  that is far from the decision boundaries should tolerate larger smoothing (and equivalently have a larger certification radius) as compared to inputs closer to these boundaries, we optimize for the amount of smoothing per input (specifically  $\sigma_x$ ) that maximizes the certification radius. This proposed process is denoted as *data dependent smoothing* where we provide a procedure for certifying the resultant smooth classifier.

<sup>2</sup>The paper mainly focuses on Gaussian smoothing, but the idea holds for other parameterized distributions.

### 3 DATA DEPENDENT SMOOTHING

#### 3.1 PRELIMINARIES AND NOTATIONS

Let  $x \in \mathbb{R}^d$  and the labels  $y \in \mathcal{Y} = \{1, \dots, k\}$  be the input-label pairs  $(x, y)$  sampled from an unknown data distribution. Unless explicitly mentioned, we consider a classifier  $f_\theta : \mathbb{R}^d \rightarrow \mathcal{P}(\mathcal{Y})$  parameterized by  $\theta$  where  $\mathcal{P}(\mathcal{Y})$  is a probability simplex over  $k$  labels. We say that  $f_\theta$  is  $\ell_p^r$  certifiably accurate for an input  $x$ , if and only if,  $\arg \max_c f_\theta^c(x) = \arg \max_c f_\theta^c(x + \delta) = y \quad \forall \|\delta\|_p \leq r$ , where  $f_\theta^c$  is the  $c^{\text{th}}$  element of  $f_\theta$ . That is to say, the classifier correctly predicts the label of  $x$  and enjoys a constant prediction for all perturbations  $\delta$  that are in the  $\ell_p$  ball of radius  $r$  from  $x$ . As such, the overall  $\ell_p^r$  certification accuracy is defined as the average certified accuracy over the data distribution. Following prior art [Cohen et al., 2019, Salman et al., 2019a, Zhai et al., 2020], we focus on  $\ell_2^r$  certification.

#### 3.2 OVERVIEW OF RANDOMIZED SMOOTHING

Randomized smoothing constructs a certifiable classifier  $g_\theta$  by smoothing a base classifier  $f_\theta$ . For any  $\sigma > 0$ , the smooth classifier is defined as:  $g_\theta(x) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)} [f_\theta(x + \epsilon)]$ . Let  $g_\theta$  predict label  $c_A$  for input  $x$  with some confidence, *i.e.*  $\mathbb{E}_\epsilon [f_\theta^{c_A}(x + \epsilon)] = p_A \geq p_B = \max_{c \neq c_A} \mathbb{E}_\epsilon [f_\theta^c(x + \epsilon)]$ , then,  $g_\theta$  is certifiably robust at  $x$  with certification radius:

$$R = \frac{\sigma}{2} (\Phi^{-1}(p_A) - \Phi^{-1}(p_B)). \quad (1)$$

Here,  $g(x + \delta) = g(x) \quad \forall \|\delta\|_2 \leq R$ , where  $\Phi$  is the CDF of the standard Gaussian.

#### 3.3 ROBUSTNESS-ACCURACY TRADE-OFF

Note that Equation 1 holds regardless of the prediction  $c_A$  made by the smooth classifier  $g_\theta$ . This suggests that one can perhaps improve the robustness of  $g_\theta$ , *i.e.* increase certification radius  $R$  where  $g_\theta$  is constant, by increasing the hyper parameter  $\sigma$  in Equation 1. However, to reason about  $\ell_2^r$  certification accuracy, it is not enough to increase the certification radius  $R$ , as this requires  $c_A$  to be the correct prediction for  $x$  by  $g_\theta$ . This reveals the robustness-accuracy trade-off as one cannot improve  $\ell_2^r$  certified accuracy by only increasing the certification radius  $R$  (robustness) through the increase in  $\sigma$ . This is because it comes at the expense of requiring a classifier  $g_\theta$  that correctly classifies  $x$  with correct label  $y$  under large Gaussian perturbations (accuracy). As such, the following inequality should hold  $\mathbb{E}_\epsilon [f_\theta^y(x + \epsilon)] \geq p_A \geq p_B \geq \max_{c \neq y} \mathbb{E}_\epsilon [f_\theta^c(x + \epsilon)]$ .

---

#### Algorithm 1: Data Dependent Certification

---

**Function** OptimizeSigma( $f_\theta, x, \alpha, \sigma_0, n$ ):

```

Initialize:  $\sigma_x^0 \leftarrow \sigma_0, K$ 
for  $k = 0 \dots K - 1$  do
    sample  $\hat{\epsilon}_1, \dots, \hat{\epsilon}_n \sim \mathcal{N}(0, I)$ 
     $\psi(\sigma_x^k) = \frac{1}{n} \sum_{i=1}^n f_\theta(x + \sigma_x^k \hat{\epsilon}_i)$ 
     $E_A(\sigma_x^k) = \max_c \psi^c; y_A = \arg \max_c \psi^c;$ 
     $E_B(\sigma_x^k) = \max_{c \neq y_A} \psi^c$ 
     $R(\sigma_x^k) = \frac{\sigma_x^k}{2} (\Phi^{-1}(E_A) - \Phi^{-1}(E_B))$ 
     $\sigma_x^{k+1} \leftarrow \sigma_x^k + \alpha \nabla_{\sigma_x^k} R(\sigma_x^k)$ 
 $\sigma_x^* \leftarrow \sigma_x^K$ 
return  $\sigma_x^*$ 

```

---

#### 3.4 DATA DEPENDENT SMOOTHING FOR CERTIFICATION

The certification region  $\mathcal{R} = \{\delta : \|\delta\|_2 \leq R\}$  at an input  $x$  is fully characterized by the classifier  $f_\theta$  and the standard deviation of the Gaussian distribution  $\sigma$ . Moreover, for a given  $f_\theta$ , the certification region  $\mathcal{R}$  varies at different  $x$ , when  $\sigma$  is fixed, due to the nonlinear dependence of the prediction gap  $\Phi^{-1}(p_A(x; \sigma)) - \Phi^{-1}(p_B(x; \sigma))$  on  $x$ . This hints that, for a given  $f_\theta$ , different inputs  $x$  may enjoy a different optimal  $\sigma_x^*$  that maximizes the certification region. To see this, consider the three inputs  $x_1, x_2$  and  $x_3$  all correctly classified by the binary classifier  $f_\theta$  as  $\mathcal{C}_1$  in Figure 1. Using a fixed  $\sigma$  to smooth the predictions of  $f_\theta$ , *i.e.* predict with  $g_\theta$ , reveals that inputs, depending on how close they are from the decision boundaries, can enjoy different levels of smoothing without affecting the prediction of  $g_\theta$ . For instance, as shown in Figure 1 for constant  $\sigma$ , the input far from the decision boundary  $x_1$  could have still been classified correctly with similarly large prediction gap even if  $f_\theta$  were to be smoothed with a larger  $\sigma$ . This indicates that perhaps the certification radius at  $x_1$  could have been enlarged with a larger smoothing  $\sigma$ . As for  $x_2$ , we can observe that while the prediction under this choice of  $\sigma$  by  $g_\theta$  is still correct, the prediction gap  $\Phi^{-1}(p_A(x; \sigma)) - \Phi^{-1}(p_B(x; \sigma))$  drops, due to having more Gaussian samples fall in the  $\mathcal{C}_2$  region. Thus, a different choice of  $\sigma$  could have been used to trade-off the drop in prediction gap and certification radius. Last, for the input  $x_3$  that is very close to the decision boundary, the sub optimal choice of  $\sigma$  (too large for  $x_3$ ) could result in an incorrect prediction by  $g_\theta$ . Despite the observations that  $\sigma$  plays a significant role in  $\ell_2^r$  certification accuracy, certification methods generally (i) choose  $\sigma$  arbitrarily and (ii) set it to be constant for all  $x$ . Based on this observation, for a given smooth classifier with a specific  $\sigma_0$ , where  $\sigma_0$  can be zero reducing the smooth classifier to  $f_\theta$ , we seek to construct another smooth classifier with parameter  $\sigma_x^*$  for every input  $x$  such that: (i) the prediction of both smooth classifiers (smoothing with  $\sigma_0$  and  $\sigma_x^*$ ) is

identical for all  $x$ . **(ii)** The certification radius of the new smooth classifier at every  $x$  is maximized. To construct a classifier smoothed with  $\sigma_x^*$  enjoying the two previous properties, let  $c_A$  be the prediction under  $\sigma_0$  smoothing, *i.e.*  $c_A = \arg \max_c \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma_0 I)} [f^c(x + \epsilon)]$ . We maximize  $R$  in Equation 1 over  $\sigma$  for every  $x$  by solving:

$$\sigma_x^* = \arg \max_{\sigma} \frac{\sigma}{2} \left( \Phi^{-1} \left( \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)} [f_{\theta}^{c_A}(x + \epsilon)] \right) - \Phi^{-1} \left( \max_{c \neq c_A} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)} [f_{\theta}^c(x + \epsilon)] \right) \right). \quad (2)$$

Since  $\Phi^{-1}$  is a strictly increasing function, it is important to note that solving Equation 2 for a fixed  $c_A$  can at worst yield a smooth classifier of an identical radius to when the classifier is smoothed with  $\sigma_0$  both predicting  $c_A$  for  $x$ .

**Solver.** While our proposed Objective 2 has a similar form to the MACER regularizer [Zhai et al., 2020] used during training, ours differs in that we optimize  $\sigma$  for every  $x$  and not the network parameters  $\theta$ , which are fixed here. A natural solver for 2 is stochastic gradient ascent with the expectation approximated with  $n$  Monte Carlo samples. As such, the gradient of the objective at the  $k^{\text{th}}$  iteration will be approximated as follows:  $\nabla_{\sigma^k} \frac{\sigma^k}{2} [\Phi^{-1}(\gamma^{c_A}(\sigma^k)) - \Phi^{-1}(\max_{c \neq c_A} \gamma^c(\sigma^k))]$ , where  $\gamma^c(\sigma^k) = \frac{1}{n} \sum_{i=1}^n f^c(x + \epsilon_i)$  for  $\epsilon_1, \dots, \epsilon_n \sim \mathcal{N}(0, (\sigma^k)^2 I)$ . However, this estimation of the gradient suffers from high variance due to the dependence of the expectation on the optimization variable  $\sigma$  that parameterizes the smoothing distribution  $\mathcal{N}(0, \sigma^2 I)$  [Williams, 1992]. To alleviate this, we use the *reparameterization trick* suggested by Kingma and Welling [2014], Rezende et al. [2014] to compute a lower variance gradient estimate for our Objective 2. In particular, with the change of variable  $\epsilon = \sigma \hat{\epsilon}$  where  $\hat{\epsilon} \sim \mathcal{N}(0, I)$ , Objective 2 is equivalent to:

$$\sigma_x^* = \arg \max_{\sigma} \frac{\sigma}{2} \left( \Phi^{-1} \left( \mathbb{E}_{\hat{\epsilon} \sim \mathcal{N}(0, I)} [f_{\theta}^{c_A}(x + \sigma \hat{\epsilon})] \right) - \Phi^{-1} \left( \max_{c \neq c_A} \mathbb{E}_{\hat{\epsilon} \sim \mathcal{N}(0, I)} [f_{\theta}^c(x + \sigma \hat{\epsilon})] \right) \right) \quad (3)$$

Note that, unlike before, the expectation over the distribution  $\hat{\epsilon} \sim \mathcal{N}(0, I)$  no longer depends on the optimization variable  $\sigma$ . This allows the gradient of 3 to enjoy a lower variance compared to the gradient of 2 [Kingma and Welling, 2014, Rezende et al., 2014]. Algorithm 1 summarizes the updates for optimizing  $\sigma$  for each  $x$  by solving 3 with  $K$  steps of stochastic gradient ascent. It is worthwhile to mention that the function `OptimizeSigma` in Algorithm 1 is agnostic of the choice of architecture  $f_{\theta}$  and of the training procedure that constructed  $f_{\theta}$ .

---

### Algorithm 2: Training with Data Dependent $\sigma_{x_i}$

---

**Function** `TrainBatch` ( $f_{\theta}$ ,  $\{x_i, y_i\}_{i=1}^B$ ,  $\{\sigma_{x_i}\}_{i=1}^B$ ,  $\alpha$ ,  $n$ ):

```

for  $i = 1, \dots, B$  do
   $\sigma_{x_i}^* = \text{OptimizeSigma}(f_{\theta}, x_i, \alpha, \sigma_{x_i}, n)$ 
  TrainFunction ( $\{x_i, y_i\}_{i=1}^B$ ,  $\{\sigma_{x_i}^*\}_{i=1}^B$ )
  // any training routine e.g. MACER

```

---

## 3.5 MEMORY-BASED CERTIFICATION FOR DATA DEPENDENT CLASSIFIERS

Unlike previous approaches where  $\sigma$  is constant for all inputs, the data dependent classifier  $g_{\theta}$  with varying  $\sigma$  per input can not be directly certified by the classical Monte Carlo algorithms proposed by Cohen et al. [2019]. This is since the data dependent classifier  $g_{\theta}$  does not enjoy a constant  $\sigma$  within the given certification region, *i.e.*  $g_{\theta}$  tailors a new  $\sigma_x$  for every input  $x$  including within the certified region of  $x$ . Informally, let  $R(\sigma_{x_1}^*)$  be the radius of certification at  $x_1$  granted by the data dependent classifier  $g_{\theta}$ . The data dependent classifier *does not guarantee* that there *can not exist*  $x_2$  within the region of certification of  $x_1$ , *i.e.*  $\|x_1 - x_2\|_2 \leq R(\sigma_{x_1}^*)$ , where  $g_{\theta}$  with  $\sigma_{x_2}^*$  predicts  $x_2$  differently from  $x_1$  breaking the soundness of certification. To circumvent this problem, we propose a memory-based procedure to certifying our proposed data dependent classifier. Let  $\{x_i\}_{i=1}^N$  be a set of previously predicted inputs and  $\{C_i\}_{i=1}^N$  be their corresponding predictions with mutually exclusive  $\ell_2$  certified regions  $\mathcal{R}_i$  for differently predicted inputs, *i.e.*  $\mathcal{R}_i \cap \mathcal{R}_j = \emptyset \forall i \neq j, C_i \neq C_j$ . Let  $x_{N+1}$  be a new input with a certified region  $\mathcal{R}_{N+1}$  computed by the Monte Carlo algorithms of Cohen et al. [2019] for the data dependent classifier  $g_{\theta}$  with prediction  $C_{N+1}$ . If there exists an  $i$  such that  $\mathcal{R}_{N+1} \cap \mathcal{R}_i \neq \emptyset$ ,  $x_{N+1} \in \mathcal{R}_i$ , and  $C_{N+1} \neq C_i$ , we adjust the prediction of the data dependent classifier  $g_{\theta}$  to be  $C_i$  and update  $\mathcal{R}_{N+1}$  to be the largest subset of  $\mathcal{R}_{N+1}$  that is a subset of  $\mathcal{R}_i$  (see middle example in Figure 2). On the other hand, if  $\mathcal{R}_{N+1} \cap \mathcal{R}_i = \emptyset$ ,  $x_{N+1} \notin \mathcal{R}_i$ , and that  $C_{N+1} \neq C_i$ , we update  $\mathcal{R}_{N+1}$  to be the largest subset of  $\mathcal{R}_{N+1}$  not intersecting with  $\mathcal{R}_i$  (see right example in Figure 2). We perform the previous operations for all elements in the memory and add  $x_{N+1}, C_{N+1}, \mathcal{R}_{N+1}$  to memory. The aforementioned procedure grants a sound certification for the data dependent classifier preventing by construction overlapping certified regions with different predictions. While the memory-based certification is essential for a sound certification, empirically, we never found in any of the later experiments a case where two inputs predicted differently suffer from intersecting certified regions. That is to say while our sound certificate works on the memory-enhanced data dependent smooth classifier, we found that the certified radius of the memory classifier for every input is the radius granted by the Monte Carlo certificates of Cohen

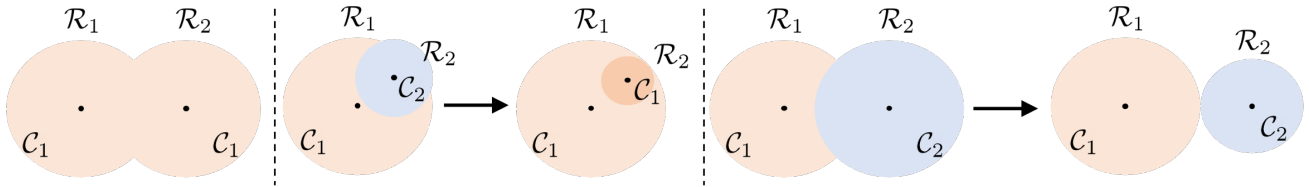


Figure 2: **Memory-based certification of the data dependent classifier.** Given a memory of an input  $x_1$  with a certified region  $\mathcal{R}_1$  and another input  $x_2$  with a certified region  $\mathcal{R}_2$ . Three scenarios could arise where  $\mathcal{R}_1$  and  $\mathcal{R}_2$  intersect. **Left:** The certified regions intersect while both  $x_1$  and  $x_2$  share the same prediction. In this case,  $x_2$  along with its certified region are directly added to memory. **Middle:**  $x_2$  lies inside  $\mathcal{R}_1$  with a different prediction from  $x_1$ . In this case,  $x_2$  is predicted with the same prediction as  $x_1$  and added to memory along with the largest subset of  $\mathcal{R}_2$  that is within  $\mathcal{R}_1$ . **Right:**  $x_2$  lies outside the  $\mathcal{R}_1$  with a different prediction from  $x_1$ . In this case,  $x_2$  with its prediction are added to memory along with the largest certified region in  $\mathcal{R}_2$  not intersecting with  $\mathcal{R}_1$ .

et al. [2019] for the data dependent classifier. Therefore and throughout, we refer to the memory-enhanced data dependent smooth classifier and data dependent smooth classifier interchangeably. We elaborate more on this and provide an algorithm in the **Appendix**.

### 3.6 TRAINING WITH DATA DEPENDENT SMOOTHING

Models that enjoy a large  $\ell_2^r$  certification accuracy under the randomized smoothing framework need to enjoy a large certification radius  $R$  in Equation 1 for all  $x$  and be able to correctly classify inputs corrupted with Gaussian noise, *i.e.*  $g_\theta(x) = y$ . While there are several approaches to train  $f_\theta$  (or directly  $g_\theta$ ) so as to output correct predictions for inputs corrupted with noise sampled from  $\mathcal{N}(0, \sigma^2 I)$ , all existing works fix  $\sigma$  for all inputs during training. We are interested in complementing these approaches with smoothing distributions that are data dependent. As such, we can employ the training procedure of these approaches but with  $\sigma_x^*$  computed by `OptimizeSigma`. Algorithm 2 summarizes this proposed training pipeline. The function `TrainFunction` proceeds by performing backpropagation using any training scheme, given the estimated  $\sigma_{x_i}^*$  for each  $x_i$ . We note that whenever Algorithm 2 is used, we initialize  $\sigma_{x_i}$  at each epoch with  $\sigma_{x_i}^*$  computed at the previous epoch. Since COHEN, SMOOTHADV and MACER are among the most popular approaches that embed randomized smoothing certificates as part of the training routine, `TrainFunction` refers here to any of these three training methods. Empirically, we show that we can boost all three methods even further when models are trained with Algorithm 2.

## 4 EXPERIMENTS

We conduct two sets of experiments to validate our key contributions. (i) We show that we can boost certified accuracy for several pre-trained models by using Algorithm 1 for data

dependent smoothing only during certification, *i.e.* without employing any additional training. (ii) Once data dependent smoothing is employed during training, we can improve the certified accuracy even further. Since our framework is agnostic to the training routine, we incorporate it into (i) COHEN [Cohen et al., 2019], (ii) SMOOTHADV [Salman et al., 2019a] and (iii) MACER [Zhai et al., 2020]. Throughout, we use DS to refer to when data dependent smoothing is used only in certification and DS<sup>2</sup> when it is used during both training and certification.

**Setup.** We conduct experiments with ResNet-18 and ReNet-50 [He et al., 2016] on CIFAR10 [Krizhevsky and Hinton, 2009] and ImageNet [Russakovsky et al., 2015], respectively. For CIFAR10 experiments, we train from scratch for 200 epochs. For ImageNet, we initialize using the network parameters provided by the authors. When  $\sigma$  is fixed and following prior art, *e.g.* COHEN, SMOOTHADV, and MACER, we set  $\sigma \in \{0.12, 0.25, 0.50\}$  and  $\sigma \in \{0.25, 0.50, 1.0\}$  for CIFAR10 and ImageNet, respectively, for training and certification. We set  $\alpha = 10^{-4}$  in Algorithm 1 and the initial  $\sigma_0$  to the  $\sigma$  used in training the respective model. Unless stated otherwise, we set  $n = 1$  in Algorithm 1. Following COHEN and SMOOTHADV, we compare models using the approximate certified accuracy curve (simply referred to as certified accuracy) followed by the envelope curve over all  $\sigma$ . We also report the Average Certified Radius (ACR) proposed by MACER  $1/|S_{test}| \sum_{(x,y) \in S_{test}} R(f_\theta, x) \cdot \mathbb{1}\{\arg \max_c g_\theta^c(x) = y\}$ , where  $\mathbb{1}\{\cdot\}$  is an indicator function. Following COHEN and all randomized smoothing methods, we certify all results using  $N_0 = 100$  Monte Carlo samples for prediction and  $N = 100,000$  estimation samples to estimate the radius with a failure probability of 0.001 given a smoothing  $\sigma$ .

### 4.1 COHEN + DS

We combine data dependent smoothing with COHEN. Following Gaussian augmentation, this method trains  $f_\theta$  on  $(x + \epsilon)$ , where  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ , with the cross entropy loss.

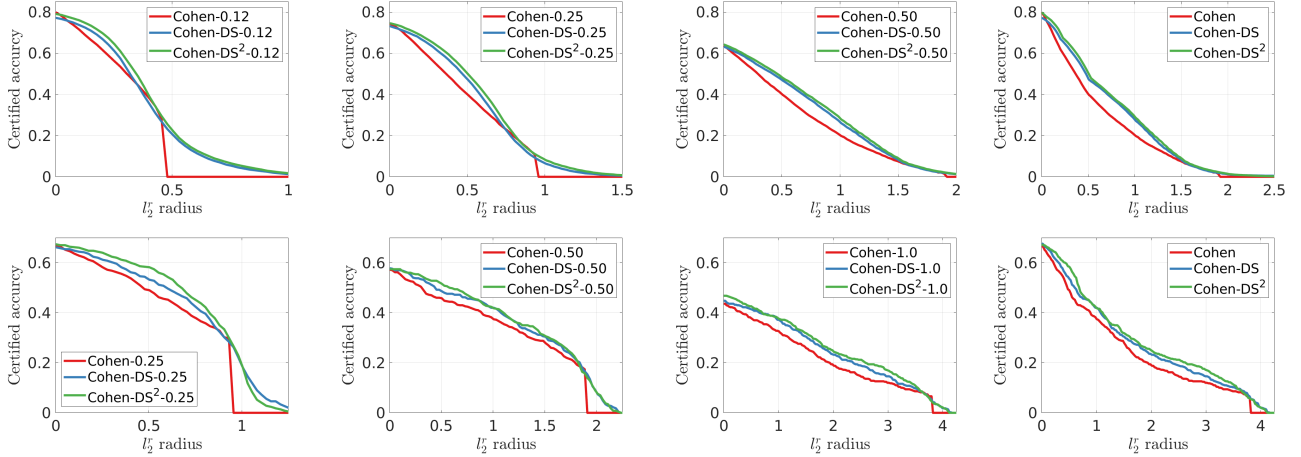


Figure 3: **Certified accuracy comparison against Cohen per radius per  $\sigma$ .** We compare Cohen against our data dependent certification Cohen-DS and when data dependency is incorporated in both training and certification Cohen-DS<sup>2</sup> for several  $\sigma$ . The value of  $\sigma$  shown for our models in the legend refers to the optimization initialization  $\sigma_0$  in Algorithm 1. We show CIFAR10 and ImageNet results in first and second rows, respectively, where the last column is the envelope.

Table 1: **Best certified accuracy per radius and ACR of Cohen, Cohen-DS and Cohen-DS<sup>2</sup>.**

CIFAR10	Radius		0.0	0.25	0.50	0.75	1.00	1.25	1.50	1.75	2.00	2.25	2.50	ACR
	Train	Certify												
Cohen	FS	FS	<b>79.9</b>	58.3	40.1	29.2	20.2	13.1	7.3	3.3	0.0	0.0	0.0	0.591
Cohen-DS	FS	DS	77.2	64.5	47.8	38.3	27.6	16.5	8.0	3.2	1.2	<b>0.7</b>	<b>0.5</b>	<b>0.784</b>
Cohen-DS <sup>2</sup>	DS	DS	79.8	<b>66.5</b>	<b>50.4</b>	<b>39.2</b>	<b>29.1</b>	<b>18.3</b>	<b>8.8</b>	<b>3.8</b>	<b>1.4</b>	0.6	0.2	0.764

ImageNet	Radius		0.0	0.25	0.50	0.75	1.00	1.50	2.0	2.5	3.0	3.50	4.0	ACR
	Train	Certify												
Cohen	FS	FS	66.6	58.2	49.0	42.4	37.4	27.8	19.4	14.4	12.0	8.6	0.0	1.098
Cohen-DS	FS	DS	<b>67.8</b>	61.4	53.6	45.6	<b>42.0</b>	30.4	23.4	18.8	14.6	10.2	<b>2.0</b>	1.257
Cohen-DS <sup>2</sup>	DS	DS	67.4	<b>64.2</b>	<b>58.4</b>	<b>47.4</b>	41.8	<b>31.8</b>	<b>25.0</b>	<b>21.2</b>	<b>17.2</b>	<b>11.0</b>	<b>2.0</b>	<b>1.319</b>

**DS for certification only.** We first certify the trained models with the same fixed  $\sigma$  used in training for all inputs, dubbed COHEN. Then, we certify using the memory based certification the same trained models with the proposed data dependent  $\sigma_x^*$  produced by Algorithm 1, which we refer to as COHEN-DS. Figure 3 plots the certified accuracy for CIFAR10 and ImageNet in the first and second rows, respectively. Even though the base classifier  $f_\theta$  is identical for COHEN and COHEN-DS, Figure 3 shows that COHEN-DS is superior to COHEN in certified accuracy across almost all radii and for all training  $\sigma$  on both datasets. This is also evident from the envelope plots in the last column of Figure 3. In Table 1, we report the best certified accuracy per radius over all training  $\sigma$  for COHEN (envelope figure) against our best COHEN-DS, cross-validated over all training  $\sigma$  and the number of iterations in Algorithm 1  $K$ , accompanied with the corresponding ACR score. For instance, we observe that data dependent certification COHEN-DS can significantly boost certified accuracy at radii 0.5 and 0.75 by 7.7% (from

40.1 to 47.8) and 9.1% (from 29.2% to 38.3%), respectively, and by 0.193 ACR points on CIFAR10. Moreover, we boost the certified accuracy on ImageNet by 4.6% and 3.2% at 0.5 and 0.75 radii, respectively, and by 0.159 ACR points.

**DS for training and certification.** We employ data dependent smoothing in both training and certification for COHEN models (denoted as COHEN-DS<sup>2</sup>) by running Algorithm 2. For CIFAR10, we train COHEN first with fixed  $\sigma$  for 50 epochs, *i.e.*  $K = 0$  in Algorithm 1, and then we perform data dependent smoothing with  $K = 1$  for the remaining 150 epochs. For ImageNet experiments, we only finetune the provided models for 30 epochs using Algorithm 2 with  $K = 1$ . Once training is complete, we certify all trained models with Algorithm 1 using the memory based certification. In Figure 3, we observe that COHEN-DS<sup>2</sup> can further improve certified accuracy across all trained models on both CIFAR10 and ImageNet. This is also evident in the last column of Figure 3 that shows the best certified accuracy per radius

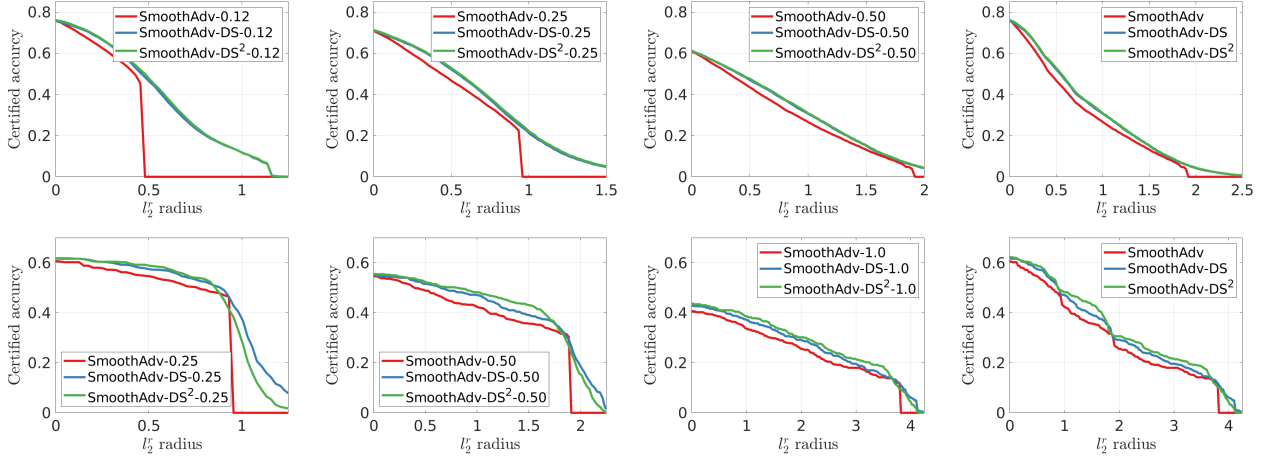


Figure 4: **Certified accuracy comparison against SmoothAdv per radius per  $\sigma$ .** We compare SmoothAdv against SmoothAdv-DS and SmoothAdv-DS<sup>2</sup>. We show CIFAR10 and ImageNet results in first and second rows, respectively.

Table 2: **Best certified accuracy per radius and ACR** of SmoothAdv, SmoothAdv-DS and SmoothAdv-DS<sup>2</sup>.

CIFAR10	Radius		0.0	0.25	0.50	0.75	1.00	1.25	1.50	1.75	2.00	2.25	2.50	ACR
	Train	Certify												
SmoothAdv	FS	FS	76.0	62.4	46.7	34.6	26.5	19.5	12.9	7.5	0.0	0.0	0.0	0.681
SmoothAdv-DS	FS	DS	75.7	66.4	52.1	38.8	30.6	22.2	15.0	8.5	4.2	1.8	0.6	0.799
SmoothAdv-DS <sup>2</sup>	DS	DS	<b>76.2</b>	<b>66.8</b>	<b>52.8</b>	<b>39.3</b>	<b>30.8</b>	<b>22.6</b>	<b>15.1</b>	<b>8.8</b>	<b>4.3</b>	<b>2.0</b>	<b>0.7</b>	<b>0.812</b>

ImageNet	Radius		0.0	0.25	0.50	0.75	1.00	1.50	2.0	2.5	3.0	3.50	4.0	ACR
	Train	Certify												
SmoothAdv	FS	FS	60.8	57.8	54.6	50.4	42.2	35.6	25.6	20.4	18.0	14.2	0.0	1.287
SmoothAdv-DS	FS	DS	62.0	60.4	57.4	53.2	47.0	39.2	29.2	23.8	19.6	15.2	<b>6.2</b>	1.445
SmoothAdv-DS <sup>2</sup>	DS	DS	<b>62.2</b>	<b>60.6</b>	<b>58.8</b>	<b>54.2</b>	<b>48.2</b>	<b>43.0</b>	<b>30.6</b>	<b>25.4</b>	<b>21.6</b>	<b>18.6</b>	4.2	<b>1.514</b>

(envelope) over all training  $\sigma$ . We note that COHEN-DS<sup>2</sup> improves the certification accuracy of COHEN-DS by 2.6% and by 0.9% at radii 0.5 and 0.75 respectively on CIFAR10, and by 4.8% and 1.8% at radii 0.5 and 0.75 respectively on ImageNet. The improvements are consistently present over a wide range of radii on both datasets. We do observe that the ACR score for COHEN-DS<sup>2</sup> on CIFAR10 marginally drops compared to COHEN-DS. We believe that this is due to the fact that some inputs that are classified correctly at the small radii have an overall larger certification radius for COHEN-DS compared to COHEN-DS<sup>2</sup> on CIFAR10. Regardless, COHEN-DS<sup>2</sup> substantially outperforms COHEN by 0.173 ACR points. As compared to COHEN-DS, COHEN-DS<sup>2</sup> improves the ACR on ImageNet from 1.257 to 1.319.

## 4.2 SMOOTHADV + DS

We combine our data dependent smoothing strategy with the more effective SMOOTHADV, which trains the smoothed classifier for every  $x$  on the adversarial example  $\hat{x}$  that max-

imizes  $-\log \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)} [f_{\theta}^y(x' + \epsilon)]$ , where  $\|x' - x\| \leq \zeta$ . For CIFAR10 experiments, we follow the training procedure of SMOOTHADV, where the adversary  $\hat{x}$  is computed with 2 PGD (proximal gradient descent) steps with  $\zeta = 0.25$  and one augmented sample to estimate the expectation. For ImageNet experiments, we use the best reported models, in terms of certified accuracy, provided by the authors, which correspond to  $\zeta = 0.5$  for  $\sigma = 0.25$  and  $\zeta = 1.0$  for  $\sigma \in \{0.5, 1.0\}$ .

**DS for certification only.** Similar to COHEN, we first certify SMOOTHADV models trained with the same fixed  $\sigma$ . Then, we certify the proposed data dependent  $\sigma_x^*$  models using the memory-based certification, which we refer to as SMOOTHADV-DS. In Figure 4, we show the certified accuracy for both CIFAR10 and ImageNet in the first and second rows, respectively. The last column shows the envelopes per radius. Even though they both share the same classifier  $f_{\theta}$ , SMOOTHADV-DS significantly improves upon SMOOTHADV over all radii and all values of  $\sigma$  in training for both CIFAR10 and ImageNet. In particular, for models trained with  $\sigma = 0.25$ , SMOOTHADV achieves a zero cer-

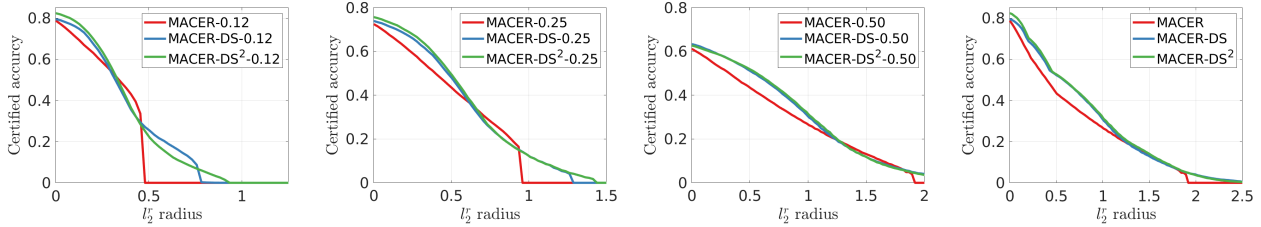


Figure 5: **Certified accuracy comparison against MACER per radius per  $\sigma$ .** We compare MACER against MACER-DS and MACER-DS<sup>2</sup> for several  $\sigma$  on CIFAR10 with the last column showing the envelope.

Table 3: **Best certified accuracy per radius and ACR** of MACER, MACER-DS and MACER-DS<sup>2</sup> on CIFAR10.

CIFAR10	Radius		0.0	0.25	0.50	0.75	1.00	1.25	1.50	1.75	2.00	2.25	2.50	ACR
	Train	Certify												
MACER	FS	FS	78.8	59.3	43.6	34.7	26.6	19.4	13.0	7.50	0.0	0.0	0.0	0.702
MACER-DS	FS	DS	79.5	66.7	52.3	43.0	30.8	19.5	12.8	7.55	<b>3.97</b>	<b>1.67</b>	<b>0.5</b>	<b>0.841</b>
MACER-DS <sup>2</sup>	DS	DS	<b>82.4</b>	<b>68.3</b>	<b>52.7</b>	<b>43.5</b>	<b>31.7</b>	<b>20.6</b>	<b>13.8</b>	<b>7.92</b>	3.65	1.39	0.4	0.807

tified accuracy for large certification radii ( $\geq 1.0$ ), while SMOOTHADV-DS achieves non-trivial certified accuracy in these cases. Similar to the earlier setup, we report the best certified accuracy along with the ACR scores in Table 2. We improve over SMOOTHADV by large margins. For example, the certified accuracy at 0.5 radius increases by 5.4% and 2.8% on CIFAR10 and Imagenet, respectively. The improvement is consistent over all radii. The ACR also improves by 0.118 and 0.158 on CIFAR10 and ImageNet, respectively.

**DS for training and certification.** We fine tune the SMOOTHADV trained models (either the retrained CIFAR10 models or the ImageNet models provided by SMOOTHADV) using Algorithm 2, where  $\sigma_x^*$  is computed using Algorithm 1. We report the per  $\sigma$  certification accuracy comparing SMOOTHADV-DS<sup>2</sup> (certified also using memory based certification) to both SMOOTHADV-DS and SMOOTHADV. SMOOTHADV-DS<sup>2</sup> further improves the certified accuracy as compared to SMOOTHADV-DS with performance gains more prominent on ImageNet. While the improvement of SMOOTHADV-DS<sup>2</sup> over SMOOTHADV-DS is indeed small, e.g. 0.7% at radius 0.5 on CIFAR10, we observe that the performance gaps are much larger on ImageNet reaching 1.4% at 0.5 radius as shown in Table 2. We see a similar trend in ACR with improvements of 0.013 and 0.069 on CIFAR10 and ImageNet, respectively. SMOOTHADV-DS<sup>2</sup> boosts the certified accuracy of SMOOTHADV at radius 0.5 by 6.1% and 4.2% on CIFAR10 and ImageNet, respectively.

### 4.3 MACER + DS

We integrate data dependent smoothing within MACER which trains  $g_\theta$  by minimizing over the parameters  $\theta$  the following objective  $-\log g_\theta(x) + \frac{\lambda}{2} \max(\gamma - \frac{2R}{\sigma}, 0) \cdot \mathbb{1}\{\arg \max_c g_\theta^c(x) = y\}$ . where  $R$

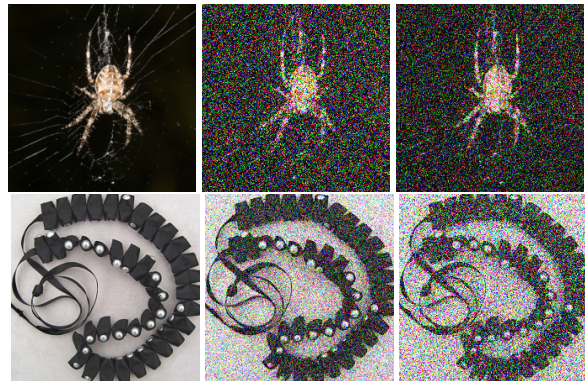


Figure 6: **Qualitative examples of estimated  $\sigma_x^*$  on different inputs.** From left to right of first row: clean image, fixed  $\sigma = 0.5$  and estimated  $\sigma_x^* = 0.368$  maximizing certification radius. Similarly for second row but with  $\sigma = 0.25$  and  $\sigma_x^* = 0.423$ . This demonstrates that  $\sigma^*$  that maximizes the radius should vary per input  $x$ .

also depends on  $\theta$ . While this seems to be similar in spirit to our approach, we in fact maximize the certification radius over  $\sigma$  with fixed parameters  $\theta$  for every  $x$ . We conduct experiments on CIFAR10 following the training procedure of MACER estimating the expectation with 64 samples,  $\lambda = 12$ , and  $\gamma = 8$ . We set  $n = 8$  in Algorithm 1 with ablations on  $n = 1$  in the **appendix**.

**DS for certification only.** Similar to the earlier setup in COHEN and SMOOTHADV, we certify models with fixed  $\sigma$  and then with data dependent  $\sigma_x^*$  using the memory based certification, referred to as MACER-DS. In Figure 5, we observe that MACER-DS significantly outperforms MACER particularly in the large radius region. This can also be seen in the envelope figure reporting the best certified accuracy



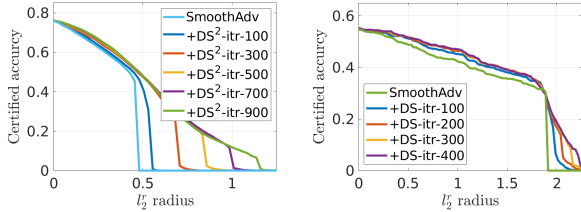


Figure 7: **Varying  $K$  in Algorithm 2.** Left figure shows certification with  $\sigma_0 = 0.12$  on CIFAR10 and  $\sigma_0 = 0.5$  on ImageNet is shown at the right.

per radius over  $\sigma$ . Similarly, Table 3 demonstrates the benefits of data dependent smoothing, where it boosts certified accuracy by 7.4% (from 59.3% to 66.7%) and 8.7% (43.6 to 52.3) at 0.25 and 0.5 radii, respectively. Moreover, we improve ACR by 0.139 points.

**DS for training and certification.** We incorporate data dependent smoothing as part of MACER training and certification in a similar fashion to the earlier setup, dubbed MACER-DS<sup>2</sup>. Figure 5 shows the improvement of MACER-DS<sup>2</sup> over the certification only MACER-DS over all trained models. Table 3 summarizes the best certified accuracy per radius. Overall, we find that the performance is comparable or slightly better than MACER-DS, which is still significantly better than MACER by 8.67% at radius 0.5. We also observe that MACER-DS enjoys better ACR than MACER-DS<sup>2</sup> with both being far better than the MACER baseline.

#### 4.4 DS FOR $\ell_1$ CERTIFICATES

At last, we extend our methodology to  $\ell_1$  certification. We leveraged the results of Yang et al. [2020] that derived the tightest  $\ell_1$  certificate using randomized smoothing with uniform distribution  $\mathcal{U}[-\lambda, \lambda]^d$ . The certified radius in that case has the form  $\mathcal{R}_1 = \lambda(p_A - p_B)$ . We replace our objective in Equation (3) with:

$$\lambda_x^* = \arg \max_{\lambda} \lambda \left( \mathbb{E}_{\epsilon \sim \mathcal{U}[-\lambda, \lambda]^d} (f_{\theta}^{c_A}(x + \epsilon)) - \max_{c \neq c_A} \mathbb{E}_{\epsilon \sim \mathcal{U}[-\lambda, \lambda]^d} (f_{\theta}^c(x + \epsilon)) \right). \quad (4)$$

We solved our objective in Eq (3) in an identical fashion to our Algorithm 1 with the same hyperparameters for  $\lambda \in \{0.25, 0.5, 1.0\}$  in certification on both CIFAR10 and ImageNet. Further, we combine our data-dependent smooth classifier with the memory based algorithm proposed in Section 3.5. It is worthwhile mentioning that similar to the  $\ell_2$  case, the memory based algorithm did not find any overlap between the certified regions of any pair of instances. We report the results in Table 4. We observe that, similar to our extensive experiments on the  $\ell_2$  certificate, our proposed

Table 4: **Best certified accuracy per  $\ell_1$  radii and ACR of YANG and YANG-DS.**

$\ell_1^r(\text{CIFAR10})$	0.0	0.25	0.5	0.75	1.0	1.5	2.0	ACR
YANG	92	83	75	71	46	0	0	0.775
YANG-DS	92	<b>89</b>	<b>82</b>	<b>76</b>	<b>58</b>	<b>6</b>	<b>2</b>	<b>0.946</b>
$\ell_1^r(\text{ImageNet})$	0.0	0.25	0.5	0.75	1.0	1.5	2.0	ACR
YANG	78	73	67	63	0	0	0	0.683
YANG-DS	<b>79</b>	<b>76</b>	<b>70</b>	<b>65</b>	<b>46</b>	0	0	<b>0.729</b>

memory-enhanced data-dependent smoothing yields consistent improvement in the  $\ell_1$  certified accuracy. We report an improvement of 7% and 3% over the state of the art certified accuracy at  $\ell_1$  radius of 0.5 on CIFAR10 and ImageNet, respectively. At last, we note similar improvement to the  $\ell_1$  ACR as reported in Table 4.

#### 4.5 DISCUSSION AND ABLATION

**Varying  $K$ .** We pose the question: does attaining better solutions to our proposed Objective 3 improve certified accuracy? To answer this, we control the solution quality of  $\sigma_x^*$  by certifying trained models with a varying number of stochastic gradient ascent iterations  $K$  in Algorithm 1. In particular, we certify the trained models SMOOTHADV-DS<sup>2</sup> and SMOOTHADV-DS on CIFAR10 and ImageNet, respectively, with a varying  $K$ . We leave the rest of the experiments for other models to the appendix. We observe in Figure 7 that the certified accuracy per radius consistently improves as  $K$  increases, particularly in the large radius regime. This is expected, since Algorithm 1 produces better optimal smoothing  $\sigma_x^*$  per input  $x$  with larger  $K$ , which in turn improves the certification radius leaving room for improvements with more powerful optimizers.

**Visualizing  $\sigma_x^*$ .** We show the variation of  $\sigma_x^*$  that maximizes the certification radius over different inputs  $x$ . Figure 6 shows two examples, where the first and fourth columns contain the clean images. In the second column, a choice of fixed  $\sigma = 0.5$  is too large compared to our estimated  $\sigma_x^* = 0.368$  that maximizes the certification radius as per Algorithm 1. As for the fifth column, we observe that a constant  $\sigma = 0.25$  is far less than  $\sigma_x^* = 0.423$ . This indicates that indeed the  $\sigma_x^*$  maximizing the certification radius varies significantly over inputs.

## 5 CONCLUSION

In this work, we presented a simple and generic framework to equip randomized smoothing techniques with data dependency. We demonstrated that combining data dependent smoothing with 3 randomized smoothing techniques provided substantial improvement in their certified accuracy.

## Acknowledgements

This publication is based upon work supported by King Abdullah University of Science and Technology (KAUST) under Award No. ORA-CRG10-2021-4648. We thank Francisco Girbal Eiras for the help in the memory based certification and the discussions.

## References

- Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *International Conference on Machine Learning (ICML)*, 2018.
- Chih-Hong Cheng, Georg Nührenberg, and Harald Ruess. Maximum resilience of artificial neural networks. In *International Symposium on Automated Technology for Verification and Analysis*, 2017.
- Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier. Parseval networks: Improving robustness to adversarial examples. *International Conference on Machine Learning (ICML)*, 2017.
- Jeremy M Cohen, Elan Rosenfeld, and J Zico Kolter. Certified adversarial robustness via randomized smoothing. *International Conference on Machine Learning (ICML)*, 2019.
- Ruediger Ehlers. Formal verification of piece-wise linear feed-forward neural networks. In *International Symposium on Automated Technology for Verification and Analysis*, 2017.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *International Conference on Learning Representations (ICLR)*, 2015.
- 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*, 2016.
- Xiaowei Huang, Marta Kwiatkowska, Sen Wang, and Min Wu. Safety verification of deep neural networks. In *International Conference on Computer Aided Verification*, 2017.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *International Conference on Learning Representations (ICLR)*, 2014.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2012.
- Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. In *IEEE Symposium on Security and Privacy (SP)*. IEEE, 2019.
- Alexander Levine and Soheil Feizi. Robustness certificates for sparse adversarial attacks by randomized ablation. In *Association for the Advancement of Artificial Intelligence (AAAI)*, 2020.
- Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive noise. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Alessio Lomuscio and Lalit Maganti. An approach to reachability analysis for feed-forward relu neural networks. *arXiv preprint arXiv:1706.07351*, 2017.
- Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *International Conference on Learning Representations (ICLR)*, 2018.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *International Conference on Machine Learning (ICML)*, 2014.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 2015.
- 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 (NeurIPS)*, 2019a.
- Hadi Salman, Greg Yang, Huan Zhang, Cho-Jui Hsieh, and Pengchuan Zhang. A convex relaxation barrier to tight robust verification of neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019b.
- Hadi Salman, Mingjie Sun, Greg Yang, Ashish Kapoor, and J Zico Kolter. Black-box smoothing: A provable defense for pretrained classifiers. *arXiv preprint arXiv:2003.01908*, 2020.

Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

Jiaye Teng, Guang-He Lee, and Yang Yuan.  $\ell_1$  adversarial robustness certificates: a randomized smoothing approach. <https://openreview.net/forum?id=H1lQIgrFDS>, 2019.

Vincent Tjeng, Kai Xiao, and Russ Tedrake. Evaluating robustness of neural networks with mixed integer programming. *International Conference on Learning Representations (ICLR)*, 2019.

Florian Tramer, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. *arXiv preprint arXiv:2002.08347*, 2020.

Jonathan Uesato, Brendan O’Donoghue, Aaron van den Oord, and Pushmeet Kohli. Adversarial risk and the dangers of evaluating against weak attacks. *International Conference on Machine Learning (ICML)*, 2018.

Tsui-Wei Weng, Huan Zhang, Hongge Chen, Zhao Song, Cho-Jui Hsieh, Duane Boning, Inderjit S Dhillon, and Luca Daniel. Towards fast computation of certified robustness for relu networks. *International Conference on Machine Learning (ICML)*, 2018.

Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 1992.

Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *International Conference on Machine Learning (ICML)*, 2018.

Greg Yang, Tony Duan, Edward Hu, Hadi Salman, Ilya Razenshteyn, and Jerry Li. Randomized smoothing of all shapes and sizes. *International Conference on Machine Learning (ICML)*, 2020.

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. *International Conference on Learning Representations (ICLR)*, 2020.

Dinghui Zhang, Mao Ye, Chengyue Gong, Zhanxing Zhu, and Qiang Liu. Filling the soap bubbles: Efficient black-box adversarial certification with non-gaussian smoothing. <https://openreview.net/forum?id=Skq8gJBFvr>, 2019.