# Promoting Robustness of Randomized Smoothing:
# Two Cost-Effective Approaches

Linbo Liu*
AWS AI Labs, UC San Diego
linbol@amazon.com

Trong Nghia Hoang
Washington State University
trongnghia.hoang@wsu.edu

Lam M. Nguyen
IBM Research
LamNguyen.MLTD@ibm.com

Tsui-Wei Weng
UC San Diego
lweng@ucsd.edu

## Abstract

*Randomized smoothing has recently attracted attentions in the field of adversarial robustness to provide provable robustness guarantees on smoothed neural network classifiers. However, existing works show that vanilla randomized smoothing usually does not provide good robustness performance and often requires (re)training techniques on the base classifier in order to boost the robustness of the resulting smoothed classifier. In this work, we propose two cost-effective approaches to boost the robustness of randomized smoothing while preserving its clean performance. The first approach introduces a new robust training method **AdvMacer** which combines adversarial training and robustness certification maximization for randomized smoothing. We show that **AdvMacer** can improve the robustness performance of randomized smoothing classifiers compared to SOTA baselines, while being 3× faster to train than MACER baseline. The second approach introduces a post-processing method **EsbRS** which greatly improves the robustness certificate based on building model ensembles. We explore different aspects of model ensembles that has not been studied by prior works and propose a novel design methodology to further improve robustness of the ensemble based on our theoretical analysis.*

## 1. Introduction

The existence of adversarial examples of deep neural networks (DNNs) [20, 7] has raised serious concerns to deploy DNNs in real-world systems, especially in the safety critical applications such as self-driving cars and aircraft control systems. Thus, many research efforts have been devoted into developing effective defenses methods to safeguard DNNs. One of the most promising direction is known as *certified defense* via *randomized smoothing*, where the word *certified* means that the defense methods have provable theoretical guarantee as opposed to easily broken heuristic defenses [3], and *randomized smoothing* is a popular technique that allows scalable certified defenses for state-of-the-art DNNs against adversarial examples. Randomized smoothing is recently proposed by [13, 15, 4] and has achieved state-of-the-art robustness guarantees. Given any classifier $f$, denoted as a *base classifier*, randomized smoothing predicts the most-likely class on the randomly perturbed input $x$ with Gaussian noises. Following this new prediction rule, randomized smoothing acts as an operator on the original *base classifier* and produce a new *smoothed* classifier which comes with provable robustness guarantees under various $\ell_p$ norm threat models [14, 4].

Unfortunately, without specially-designed training techniques, the robustness certificate of the smoothed classifier is usually very weak [4]. Thus, a few recent works [19, 22] have proposed specialized robust training algorithms to improve robustness of the smoothed classifier. In [19], the authors propose an adversarial training method called SmoothAdv, which is similar to the PGD training [17] but on the *smoothed* classifier. On the other hand, [22] propose MACER, whose training objective involves a term to maximize the robustness certificate directly. However, SmoothAdv often requires heavy tuning on a number of hyper-parameters for different noise level $\sigma$, which could be computationally challenging; while MACER usually requires longer (3×) training epochs to train and unfortunately the resulting models often have weaker certificate despite higher clean accuracy.

Motivated by the need of cost-effective robust training methods for randomized smoothing, we propose two approaches to address the limitations of existing robust training
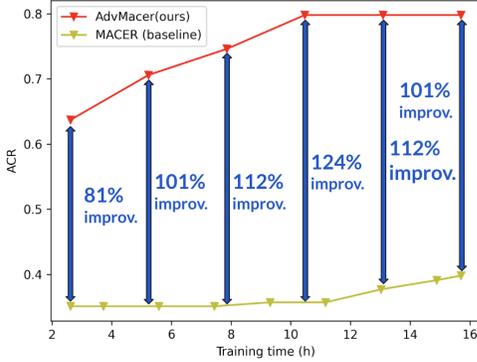
---

Figure 1: Best ACR seen so far: $x$-axis is the training time in hours and $y$-axis is the best ACR obtained before this specific training time. ACR is recorded for every 25 training epochs. Our proposed method is plotted in red, while the baseline is in light yellow. The experiments are on Cifar-10 with $\sigma = 1.00$. Improvement over baseline model is indicated in blue.

algorithms. Our contributions are three-fold: **1)** we propose a new robust training method called **AdvMacer** , which takes the best of both worlds in SmoothAdv and MACER: **AdvMacer** can achieve the best ACR while having the same computational cost as SmoothAdv and much ($3\times$) faster than MACER. If compared with MACER under same training time, our **AdvMacer** shows remarkable improvement (up to 124%) in ACR over MACER, which is illustrated in Figure 1. **2)** we equip our **AdvMacer** models with a training-free ensemble method **EsbRs**, which can further enlarge the resulting model's certified radius (by up to 8% compared with SmoothAdv and 15% compared with MACER). Crucially, we present a general theoretical analysis on ensembles and demonstrate the effect of both *intra*-model ensembles and *mixed*-model ensembles from the theoretical point of view. **3)** grounded by our theoretical findings, an optimal weighted ensemble can be derived analytically where the weights are dependent on the input data.

## 2. Related works and backgrounds

In this section, we first give backgrounds on randomized smoothing and the related certified defense SmoothAdv [19] and MACER [22]. Next, we review recent literature on applying ensemble methods to randomized smoothing.

**Randomized smoothing.** Consider a neural network classifier $f : \mathbb{R}^d \to \mathcal{Y}$ that maps an input sample $x \in \mathbb{R}^d$ to its predicted label in $\mathcal{Y}$. [4] introduced a randomized smoothing (RS) technique that can turn any base classifier $f(x)$ into a smoothed classifier $g(x)$ with provable robustness

guarantees. When taking a sample $x$, the smoothed classifier $g$ returns the class that the base classifier $f$ is most likely to return under isotropic Gaussian noise perturbation of $x$: $g(x) = \arg\max_{c \in \mathcal{Y}} \mathbb{P}_{\epsilon \sim \mathcal{N}(0,\sigma^2 I)}(f(x + \epsilon) = c)$, where $\sigma$ is the noise level that controls the trade-off between clean accuracy and model robustness. [4] further proved the robustness guarantees of such smoothed classifier in Theorem 2.1. Let $\Phi$ denote the cumulative density function (CDF) of the standard Gaussian distribution. Suppose that under Gaussian perturbation $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$, the most likely class $c_A$ is returned with probability $p_A$ and the second most likely (runner-up) class $c_B$ is returned with probability $p_B$, i.e. $c_A = \arg\max_{c \in \mathcal{Y}} \mathbb{P}(f(x + \epsilon) = c), c_B = \arg\max_{c \neq c_A} \mathbb{P}(f(x+\epsilon) = c), p_A = \mathbb{P}(f(x+\epsilon) = c_A), p_B = \mathbb{P}(f(x + \epsilon) = c_B)$.

**Theorem 2.1** (Theorem 1 of [4])**.** *Assume $p_A$ attains a lower bound $\underline{p}_A$ and $p_B$ attains an upper bound $\bar{p}_B$ with $\underline{p}_A > \bar{p}_B$, then $g(x + \delta) = c_A$ for all $\|\delta\|_2 < R$, where $R = \frac{\sigma}{2}(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\bar{p}_B))$.*

In practice, Monte Carlo sampling is employed to obtain an estimate of $p_A$, see [4]. Unfortunately, as reported in [4], the robustness certificate is weak without any specifically-designed training techniques for randomized smoothing. To enhance the robustness of randomized smoothing, [19] proposed to train base classifier $f$ on adversarial examples of soft-RS classifiers $G(x)$, which are generated by PGD [17]. Another line of work [22] considered an attack-free robust training by directly maximizing certified radius of each training sample. We briefly revisit these two methods [19, 22] in the following: Formally, let $F^\beta : \mathbb{R}^d \to P(\mathcal{Y})$ be the soft version of classifier $f$ whose last layer is a softmax layer with inverse temperature $\beta$ and $P(\mathcal{Y})$ is a probability distribution over the label space $\mathcal{Y}$. We omit the superscript $\beta$ if there is no ambiguity and denote a smoothed soft classifier as $G(x) = \mathbb{E}_{\delta \sim \mathcal{N}(0,\sigma^2 I)} F(x + \delta)$.

**SmoothAdv.** [19] introduced SmoothAdv to find adversarial examples by PGD. Denote $L_{\text{CE}}$ as the canonical cross entropy loss. Given a labeled data $(x, y)$, SmoothAdv finds a point $\hat{x}$ that maximizes the cross entropy loss of $G(x)$ in the local neighborhood of $x$:

$$\hat{x} = \arg\max_{\|x'-x\|_2 \leq \epsilon} L_{\text{CE}}(G(x'), y)$$
$$= \arg\max_{\|x'-x\|_2 \leq \epsilon} -\log \mathbb{E}_{\delta \sim \mathcal{N}(0,\sigma^2 I)} F(x' + \delta)_y. \quad (1)$$

Such optimization problem (1) is solved by projected gradient descent (PGD). To estimate the gradient of (1), [19] used Monte Carlo simulation to approximate $\nabla_{x'} L_{\text{CE}}(G(x'), y)$ by $\nabla_{x'}\left(-\log\left(\frac{1}{m}\sum_{k=1}^m F(x' + \delta_k)_y\right)\right)$, where $\delta_1, \dots, \delta_m$ are drawn i.i.d. from $\mathcal{N}(0, \sigma^2 I)$.
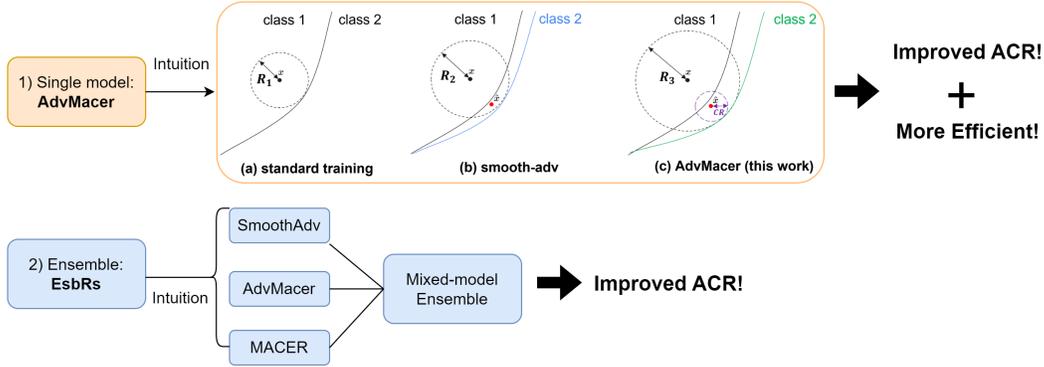
Figure 2: The overview figure, including the illustration of the idea behind **AdvMacer** : $x$ (black dot) is the original data sample and $\hat{x}$ (red dot) is an adversarial example of $x$. The solid black line is the original decision boundary. The blue line in (b) is the decision boundary using SmoothAdv and the green line in (c) is the decision boundary after applying **AdvMacer** . SmoothAdv force the classifier to classify $\hat{x}$ correctly to get the red boundary. **AdvMacer** force $\hat{x}$ to not only have correct prediction but also a large margin. Therefore, **AdvMacer** can obtain larger certified radius $R_3 >$ certified radius of smoothadv $R_2 >$ certified radius of the original classifier $R_1$.

**MACER.** Since the certified radius is related to the difference between the top probability $p_A$ and the runner-up probability $p_B$, [22] constructed MACER loss $L_{\text{MACER}}$, which aims at simultaneously minimizing classification error and maximizing the certified radius of correctly classified samples. Specifically,

$$L_{\text{MACER}}(x) = L_{\text{CE}}(G(x), y) + \lambda L_{\text{R}}(G; x, y), \quad (2)$$

where $\lambda \geq 0$ is a tuning parameter. The loss in (2) involves the soft smoothed classifier $G$ and [22] proposes to approximate $G(x)$ by Monte Carlo sampling:

$$G(x) \approx \hat{z}(x) = \frac{1}{m} \sum_{k=1}^{m} F(x + \delta_k),$$

$$\hat{g}(x) = \arg\max_{i \in \mathcal{Y}} \hat{z}_i(x). \quad (3)$$

where $\delta_1, \ldots, \delta_m$ are drawn i.i.d. from $\mathcal{N}(0, \sigma^2 I)$. Denote $\widehat{\text{CR}}(x, y)$ as the approximated certified radius at $x$, and $\widehat{\text{CR}}(x, y) = \frac{\sigma}{2}(\Phi^{-1}(\hat{z}_y(x)) - \Phi^{-1}(\max_{y' \neq y} \hat{z}_{y'}(x)))$. Therefore, the robustness loss $L_{\text{R}}(G; x, y)$ can also be approximated by

$$L_{\text{R}}(\hat{z}; x, y) = \phi(\epsilon + \tilde{\epsilon} - \widehat{\text{CR}}(x, y)) \mathbf{1}(\hat{g}(x) = y)$$

$$= \frac{\sigma}{2}\phi(\gamma - \hat{\xi}_\theta(x, y)) \mathbf{1}(\hat{g}(x) = y), \quad (4)$$

where $\epsilon, \tilde{\epsilon} > 0$ are hyper-parameters in the hinge loss, $\phi(u) = \max\{u, 0\}$, $\gamma = \frac{2}{\sigma}(\epsilon + \tilde{\epsilon})$, and $\hat{\xi}_\theta(x, y) = \Phi^{-1}(\hat{z}_y(x)) - \Phi^{-1}(\max_{y' \neq y} \hat{z}_{y'}(x))$. Finally, MACER trains a base classifier by minimizing the approximated MACER loss on training dataset. We refer readers to [22] for more details.

**Other related works.** [11] introduced consistency loss as a regularization to improve the robustness of RS classifiers. Also, the usage of mixture of adversarial examples and clean examples is suggested in training to increase robust certification as in [10]. As can be shown in Sec. 4, our proposed **AdvMacer** outperforms all baselines on Cifar-10.

**Boosting robustness via ensemble.** Model ensemble is a popular technique in the machine learning literature to practically improve model performance and reduce generalization errors [2]. Recently, there are a few works investigating the idea of using model ensemble to improve robustness of a randomized smoothed classifier [9, 21]. However, [9] mainly focused on ensemble the same type of models (i.e. models trained from the same process but with different random seeds) and only gave a brief exploration (in their App. G3.5) of mixed-model ensemble. In contrast, as will be introduced in Section 3.2, our proposed **EsbRs** is a more general ensemble method where we theoretically analyze the effect of *mixed*-model ensembles. Although weighted ensemble has also been studied in [9], their model learns the weights from training and cannot justify the weights' optimality. However, in our work, we develop a novel design framework of the optimal weight ensemble by solving an optimization problem based on our theory.

## 3. Our proposed main methods

In this section, we propose two novel and cost-effective approaches to improve robustness of a randomized smoothed classifier. First, we introduce a new robust training method **AdvMacer** that aim to maximize the certified radius over

adversarial examples, while being $3\times$ faster to train than MACER due to faster convergence. We present the intuitions, formulations as well as the details of our algorithm in Sec. 3.1. Next, in Sec. 3.2, we propose a novel ensemble method called **EsbRs** with theoretical analysis. Different from the two recent works [21, 9], we provide a more general analysis which does not require individual classifiers to come from the same training method. Our analysis allows the derivation of the optimal weight for individual classifiers, which is the key to promote robustness and the study of optimal weight has not been explored in the prior work.

## 3.1. Approach 1: AdvMacer

Inspired by the prior work SmoothAdv [19] and MACER [22] and to address their limitations, we argue that a smoothed classifier can be trained to have larger certified radius by directly optimizing the certified radius of adversarial examples instead of the clean data points. Notice that this statement requires adversarial example to be predicted correctly, which is encoded in our $L_{CE}$ term (without this constraint, the certified radius of original data point may actually decrease). The intuition is illustrated in Figure 2. Based on the above idea, we propose the following formulation.

**Formulation.** Given data $x$ and its label $y$, we aim to minimize the proposed **AdvMacer** loss consisting of two terms: $L_{AdvMacer}(x) = L_{CE}(\hat{z}(\hat{x}), y) + \lambda L_R(\hat{z}; \hat{x}, y)$, where $\hat{z}$ and $L_R$ are given in Equation (3) and Equation (4) respectively. The 1st term $L_{CE}(\hat{z}(\hat{x}), y)$ is to encourage adversarial examples $\hat{x}$ to be classified correctly, and the 2nd term $L_R(\hat{z}; \hat{x}, y) = \frac{\sigma}{2} \max\{\gamma - \hat{\xi}_\theta(\hat{x}, y), 0\} \mathbf{1}(\hat{g}(\hat{x}) = y)$ is to maximize the certified radius at the adversarial example $\hat{x}$, where $\hat{x} = \arg\max_{\|x'-x\|_2 \leq \epsilon} L_{CE}(\hat{z}(x'), y)$. To minimize the $L_{AdvMacer}(x)$, we generate the adversarial examples $\hat{x}$ via Equation (1) with $T$-step PGD using SmoothAdv [19], i.e. in the $i$-th step, we update $x_{i+1} = \prod_{\mathcal{B}(x,\epsilon)}(x_i + \nabla_x(-\log(\frac{1}{m}\sum_{k=1}^{m} F(x + \delta_k)_y))|_{x=x_i})$, where $\prod_{\mathcal{S}}(\cdot)$ is the projection onto set $\mathcal{S}$ and we set $\hat{x} = x_T$. The training objective is to minimize $L_{AdvMacer}(x)$ by first-order optimization method, and a detailed algorithm is presented in Appendix A due to page constraint.

**Hyper-parameters.** Note that there are a few hyper-parameters in **AdvMacer** : $\sigma$ is the noise level that is introduced when $f$ or $F$ is smoothed; $\epsilon$ in Equation (1) controls the size of the $\ell_2$ ball when doing PGD; $\gamma$ in Equation (4) is the parameter in hinge loss; $\lambda$ is the regularization parameter which controls the trade-off between clean accuracy and robustness; $m$ in Equation (3) is the number of Monte Carlo samples used to estimate $G(x)$; $T$ is the number of PGD step to generate adversarial samples. Finally, recall that the soft classifier $F = F^\beta$, where $\beta$ is the inverse temperature in softmax layer. The larger $\beta$ is, the closer the soft classifier $F$ is to the hard classifier $f$.

**Discussion and Comparison.** **(I).** Our **AdvMacer** is NOT just a naive combination of existing work. Instead, this clever observation has intuition (as shown in Figure 2) and can achieve better results (as shown in Table 1). Our proposed **AdvMacer** trains a model on adversarial examples while taking certified radius into consideration, which bridged between robust training and adversarial training. **(II).** Compared with SmoothAdv, **AdvMacer** doesn't bring any additional computational overhead to calculate robust loss as there exist analytic formula for certified radius; in the meantime, compared with MACER, we require much fewer number of epochs ($3\times$ smaller) and faster training time (2-4 $\times$ faster) to obtain a robust model with much larger certified radius. From the experiments in Sec. 4, it can be seen that **AdvMacer** outperforms both SmoothAdv and MACER on various dataset, such as Cifar-10, ImageNet and SVHN. **(III).** Equipped with our ensemble method **EsbRs** presented in Sec. 3.2, **AdvMacer** also enriches the diversity of component models, making mixed-model ensemble more robust (see Figure 3). For a thorough comparison by experiments, see Sec. 4 for more details.

## 3.2. Approach 2: EsbRs

### 3.2.1 Analysis

Ensemble is a cost-effective post-training technique to enhance model performance and reduce generalization error without spending much additional efforts on re-training the neural networks. By simply averaging the output from several models, ensemble shows remarkable boost in test accuracy and model robustness. Recently, there are a few works investigating the idea of using model ensemble to improve robustness of a randomized smoothed classifier [9, 21]. However, the existing work mainly focused on ensembling similar classifiers (**intra-model ensemble**) with naive averaged weights. In contrast, we also consider **mixed-model ensemble** with component classifiers coming from different training methods and conduct theoretical analysis explaining the success of mixed ensemble in certain cases. Besides, unlike [16] learning the ensemble weights empirically from training set, we develop a novel theoretical framework to design optimal ensemble weights based on our analysis. Empirical experiments verify the superiority of our proposed methods.

**Formulation.** Suppose we have $k$ trained soft classifiers $F^1, \ldots, F^k : \mathbb{R}^d \to P(\mathcal{Y})$ and $\mathcal{Y} = \{1, \ldots, c\}$. Consider soft-ensemble model $H$ whose output is a weighted average of the probabilities from $F^1, \ldots, F^k$: $H(x) = \sum_{l=1}^{k} w_l(x) F^l(x)$. Suppose the associated hard classifier is $h(x) = \arg\max_{c \in \mathcal{Y}} (H(x))_c$. Note that the weights here can be either data-dependent or data-independent. When it comes to data-dependent randomized smoothing, another line of work [1, 6] consider the noise level $\sigma$ being data-

dependent, which is different from this work. A more closely related work [21] proposed Max-Margin-Ensemble (MME) that has data-dependent weight function. A comparison of our method to MME is summarized in Tab. 3. Then we apply RS to $h$ and get the corresponding smoothed classifier $e$. It can be shown that the RS classifier $e$ of an ensemble model also has the same certification guarantee as in [4], which is presented in the following theorem.

**Theorem 3.1** (Robustness gaurantee for ensemble). *Suppose that under Gaussian perturbation $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$, the most likely class $c_A$ is returned by the **EsbRs** $e$ with probability $p_A$ and the second most likely class $c_B$ is returned with probability $p_B$, then we have $e(x+\delta) = e(x)$ for all $\|\delta\|_2 \leq R$, where $R = \frac{\sigma}{2}(\Phi^{-1}(p_A) - \Phi^{-1}(p_B))$.*

The proof is given in Appendix E. Extensive experiments from Sec. 4 show that ensemble-RS (**EsbRs**) classifier $g$ noticeably improves both accuracy and robustness, no matter $F^l$ comes from the same or different training methods. Specifically, if $F^l$ comes from more than one training methods, we call $g$ a **mixed-model ensemble**.

**Theoretical analysis.** We present some theoretical analysis on how mixed-model ensemble can reduce the variance and hence increase certified radius. We generalize the analysis in [9] to allow mixed ensemble, which provide deeper insights on model ensemble study.

For a fixed query point $x$ with a Gaussian perturbation $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$, suppose probability logits vector $\boldsymbol{y}^l \in P(\mathcal{Y})$ is returned by $F^l$. Without loss of generality, assume 1 is the majority class in RS for the ensemble model $h(x)$. For simplicity, we can work with classification margin $z_i^l = y_1^l - y_i^l$, for $i \in \mathcal{Y}$. Let $\bar{\boldsymbol{y}} = H(x + \epsilon)$. Therefore, $\bar{\boldsymbol{y}} = \sum_{l=1}^k w_l \boldsymbol{y}^l$. Similarly define $\bar{z}_i = \bar{y}_1 - \bar{y}_i$. Consider $\mathbb{E}[\bar{\boldsymbol{z}}] \in \mathbb{R}^c$ and $\mathrm{Var}(\bar{\boldsymbol{z}}) \in \mathbb{R}^{c \times c}$, where the expectation is taken over the randomness in training process, including random initialization and stochasticity in GD. Then we have

$$\mathrm{Var}(\bar{\boldsymbol{z}}) = \mathrm{Var}\left(\sum_{l=1}^k w_l \boldsymbol{z}^l\right)$$
$$= \sum_{l=1}^k w_l^2 \mathrm{Var}(\boldsymbol{z}^l) + 2\sum_{l \neq m} w_l w_m \mathrm{Cov}(\boldsymbol{z}^l, \boldsymbol{z}^m) \quad (5)$$

Hence, $\mathrm{Var}(\bar{z}_i) = \mathrm{Var}(\bar{\boldsymbol{z}})_{ii}$. Denote $p_i(w) = \mathrm{Var}(\bar{z}_i)$ as a function of $w = [w_1, \ldots, w_k]^\top$. Suppose there are a fixed number of training methods and denote this number by $s$, so $\alpha_i = \alpha_i(s) = \max_{1 \leq l \leq k} \mathrm{Var}(\boldsymbol{z}^l)_{ii}, \beta_i = \beta_i(s) = \max_{l \neq m} \mathrm{Cov}(\boldsymbol{z}^l, \boldsymbol{z}^m)_{ii}$ are functions of $s$ instead. As a result, $\alpha_i(s), \beta_i(s) = O(1)$ even as $k \to \infty$.

**A special case.** As a special case, consider $w_l = \frac{1}{k}$ for all $l = 1, \ldots, k$. By Equation (5), we derive

$$p_i(w) = \mathrm{Var}(\bar{z}_i) \leq \frac{k\alpha_i + k(k-1)\beta_i}{k^2} = \beta_i + \frac{\alpha_i - \beta_i}{k}. \quad (6)$$

These classifiers either come from different training methods, or same training method with different random seeds. Thus, existing work all assumes that the logits from one classifier have larger covariance $\alpha_i$ than the logits from different classifiers $\beta_i$. However, as we will see in **Discussion** paragraph, ensemble may harm the performance if the above assumption doesn't hold. For now, let's assume $\alpha_i > \beta_i$. By Equation (6), we conclude that the upper bound of $\mathrm{Var}(z_i)$ decreases to a constant $\beta_i$ as $k \to \infty$.

Next, we explain how $\mathrm{Var}(\bar{z}_i)$ affects certified radius. From Theorem 2.1, we see that $R = \sigma \Phi^{-1}(\underline{p}_A)$ if $\underline{p}_A \geq \frac{1}{2}$, hence we only need to show a lower bound on the top class probability $p_A$ increases as $k$ becomes larger. Since we assume the majority class's number is 1, we see that

$p_1 = \mathbb{P}(\bar{z}_i > 0, \forall i = 2, \ldots, c) \geq 1 - \sum_{i=2}^c \mathbb{P}(\bar{z}_i \leq 0)$.

By Chebyshev's inequality, $\mathbb{P}(\bar{z}_i \leq 0) \leq \mathbb{P}\left(\left|\bar{z}_i - \mathbb{E}[\bar{z}_i]\right| \geq \mathbb{E}[\bar{z}_i]\right) \leq \frac{\mathrm{Var}(\bar{z}_i)}{\mathbb{E}[\bar{z}_i]^2}$ and let $e_i = e_i(s) = \min_l \mathbb{E}[z_i^l]$, thus we have

$$p_1 \geq 1 - \sum_{i=2}^c \frac{\mathrm{Var}(\bar{z}_i)}{e_i^2}. \quad (7)$$

The above equation suggests us to choose the weight $w$ that maximizes the RHS of Equation (7) to have a larger $p_1$, hence larger certified radius. Since $e_i$ is independent of the choice of $w$, we can obtain the optimal weight by solving

$$\min_{w \in \mathbb{R}^k} \sum_{i=2}^c a_i p_i(w) \quad \text{s.t.} \sum_{l=1}^k w_l = 1, \quad w_l \geq 0, \quad (8)$$

where $a_i = e_i^{-2}$ are constants. Note that when $w_l = \frac{1}{k}$ for all $l = 1, \ldots, k$, we have a lower bound on $p_1$ by (6) and (7):

$$p_1 \geq 1 - \sum_{i=2}^c \frac{\beta_i + (\alpha_i - \beta_i)/k}{e_i^2} \to 1 - \sum_{i=2}^c \frac{\beta_i}{e_i^2} \quad \text{as } k \to \infty.$$

This explains why larger $k$ makes $p_1$ and certified radius larger even in average ensemble.

**Discussion.** Compared with [9], we generalize their analysis to allow mixed-model ensemble and hence have several new findings. **(I).** If $\alpha_i < \beta_i$, namely the logits from one model have smaller variance than those from different models, the RHS of Equation (6) becomes an increasing function in $k$, which implies ensemble does not always work. **(II).** We are the first to provide **both** theoretical analysis and comprehensive investigation on the advantage of a mixed-model

ensemble: Suppose $F^1, F^2$ come from model category 1 (for example, SmoothAdv) and $F^3$ comes from model category 2 (for example, **AdvMacer** ). If the logits from different types of models have smaller variance than those from the same type of model, namely $\mathrm{Cov}(F^1, F^3) < \mathrm{Cov}(F^1, F^2)$, $\beta_i$ will become smaller and makes mixed-model ensemble work better than intra-model ensemble. This phenomenon is observed in Figure 3. Figure 3 shows potential benefits that could be brought by heterogeneity of ensemble models, which is in accordance with our theoretical analysis presented in this section. **(III).** [9] only gave a brief exploration (in their App. G3.5) suggesting similar ensemble would outperform mixed-model ensemble and did not further analyze the reasons. Our work pointed out that it's NOT always the case: we had in fact established concrete scenarios where **mixed-model** ensemble would have **benefit** over a intra-model ensemble both theoretically and empirically (see below Figure 3). It serves as a counter example to the statement in Appendix G3.5 of [9]. Importantly, we note that our theoretical analysis is *more general* than that presented in [9] and would reduce to theirs if setting $\mathrm{cov}(z^l, z^m) = \mathrm{var}(z^l)$ in Equation (5).

**Designing optimal weighted ensemble.** The optimization problem in (8) allows us to design an optimal weight that can maximize the lower bound on $p_1$. Consider the case where $k = 2$, then (8) can be solved analytically given the knowledge of $\mathrm{Var}(z^1), \mathrm{Var}(z^2), \mathrm{Cov}(z^1, z^2)$ and $a_i$. To see this, let $b_i = \mathrm{Var}(z^1)_{ii}, c_i = 2\mathrm{Cov}(z^1, z^2)_{ii}, d_i = \mathrm{Var}(z^2)_{ii}$, then the objective function in (8) can be re-written as

$$q(w) = \sum_{i=2}^{c} a_i(b_i w_1^2 + c_i w_1 w_2 + d_i w_2^2)$$
$$\overset{(i)}{=} \sum_{i=2}^{c} a_i[b_i w_1^2 + c_i w_1(1 - w_1) + d_i(1 - w_1)^2], \quad (9)$$

where (i) uses the constraint $w_1 + w_2 = 1$ to eliminate $w_2$. Therefore, the problem (8) can be further cast as a quadratic optimization with linear constraints:

$$\min_{w_1 \in \mathbb{R}} \quad q(w_1) = A w_1^2 + B w_1 + C$$
$$\text{s.t.} \quad 0 \le w_1 \le 1, \quad (10)$$

where $A = \sum_{i=2}^{c} a_i(b_i + c_i + d_i), B = \sum_{i=2}^{c} -a_i(c_i + 2d_i)$ and $C = \sum_{i=2}^{c} a_i d_i$. Notice that this problem has an analytical solution: if $A > 0$ and $0 \le -\frac{B}{2A} \le 1$, $w_1 = -\frac{B}{2A}$ and $w_2 = 1 + \frac{B}{2A}$; else $q(w_1)$ attains minimum at boundary $w_1 = 0$ or 1.

Next, we aim at giving an estimate of $a_i, b_i, c_i, d_i$. To account for randomness both from training and Gaussian perturbation $\epsilon$ around the input $x$, we first generate $n$ i.i.d. Gaussian noisy data $x_1, \ldots, x_n$ from $\mathcal{N}(x, \sigma^2 I)$. Second,
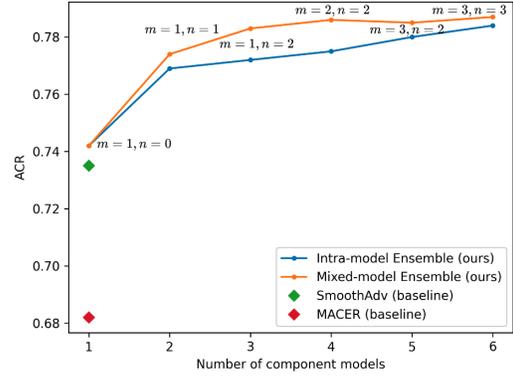


Figure 3: The plot of ACR against different number of component models in **EsbRs** on Cifar-10 with $\sigma = 0.50$. Single ensemble uses $N$ **AdvMacer** models. Mixed ensemble with totally $N$ component models uses $m$ **AdvMacer** models and $n$ SmoothAdv models.

we incorporate random perturbation for the parameters $\theta$ in classifier $F$ to imitate random seeds in training, as this is the cheapest way (without extra training cost). We randomly select $t\%$ parameters from $F$ and add i.i.d. Gaussian noise $\delta \sim \mathcal{N}(0, \tilde{\sigma}^2)$ for each selected parameter. This returns a perturbed model $\hat{F}$ from the base model $F$. Repeating the above process on $F^1$ and $F^2$ for $m$ times gives us $2m$ perturbed models $\hat{F}_1^1, \ldots, \hat{F}_m^1$ and $\hat{F}_1^2, \ldots, F_m^2$.

Now, we pass $x_1, \ldots, x_n$ into $\hat{F}_1^1, \ldots, \hat{F}_m^1$ to get $mn$ output logits vector $y^{1,1}, y^{1,2}, \ldots, y^{1,mn}$. Also, $y^{2,1}, y^{2,2}, \ldots, y^{2,mn}$ can be obtained similarly by passing $n$ noisy data into $m$ perturbed models of $F^2$. Compute $z_i^{l,j} = y_1^{l,j} - y_i^{l,j}$ for $1 \le i \le c$ and $l = 1, 2$. Then an estimation of variance and covariance can be their empirical parallel:

$$b_i = \mathrm{Var}(z^1)_{ii} = \frac{1}{mn} \sum_{j=1}^{mn} (z^{1,j} - \bar{z}^1)(z^{1,j} - \bar{z}^1)_{ii}^{\top},$$

$$c_i = 2\mathrm{Cov}(z^1, z^2)_{ii} = \frac{2}{mn} \sum_{j=1}^{mn} (z^{1,j} - \bar{z}^1)(z^{2,j} - \bar{z}^2)_{ii}^{\top},$$

$$d_i = \mathrm{Var}(z^2)_{ii} = \frac{1}{mn} \sum_{j=1}^{mn} (z^{2,j} - \bar{z}^2)(z^{2,j} - \bar{z}^2)_{ii}^{\top},$$

where $\bar{z}^l = \frac{1}{mn} \sum_{j=1}^{mn} z^{l,j}$ for $l = 1, 2$. Also obtain $a_i = e_i^{-2} = \min\{\bar{z}_i^1, \bar{z}_i^2\}^{-2}$ Hence, we can solve (10) by plugging in $a_i, b_i, c_i, d_i$. A detailed algorithm is given in Algorithm 2 in Appendix B.

**Remark 3.2.** *(I). To our best knowledge, we are the first work to develop a practical and theoretical grounded*

*methodology to obtain the optimal weight of the ensemble scheme by solving a optimization problem. We note that the two recent works [21, 9] did not explore this direction. (II). Our design strategy can be easily generalized to $k > 2$. However, due to non-convexity of the objective function in Equation (8), solving the optimization problem (8) is always intractable. Thus we only provide empirical experiments on $k = 2$ and will leave $k > 2$ as a future work.*

## 4. Experiments

In this section, we present experimental results that empirically evaluate the performance of our proposed methods, **AdvMacer** and **EsbRs**, on Cifar-10 [12], ImageNet [5] and SVHN [18] dataset. To make fair comparisons with previous baseline models, we use the same architectures as in [4]: ResNet-110 [8]. We train our models with $\sigma = 0.25, 0.50, 1.00$ on Cifar-10 and ImageNet, and $\sigma = 0.25, 0.50$ on SVHN. We train all models on a single NVIDIA V100 GPU and the training time reported below is all from NVIDIA V100 GPU. We compare the performance of **AdvMacer** with four baseline models (SmoothAdv, MACER, Consistency, SmoothMix) in Tab. 1 and observe that our **AdvMacer** achieves the largest ACR across all $\sigma$'s.

**Evaluation.** We mainly evaluate model performance on two metrics: clean accuracy and average certified radius (ACR). Clean accuracy is the classification accuracy when taking the original test images as the input and cannot evaluate model robustness. A more reasonable metric for evaluating robustness is ACR. We follow the standard evaluation protocol used in [4, 19, 22] for fair comparison: for each test data $(x_i, y_i) \in \mathcal{S}_{\text{test}}$, record the radius $R_i$ that can be certified by the model $g$. Set $R_i = 0$ if $x_i$ can't be classified correctly by $g$. Then ACR $= \frac{1}{|\mathcal{S}_{\text{test}}|} \sum_i R_i$. Since the denominator is the size of the full test set, one cannot obtain large ACR without high accuracy. Thus ACR becomes a popular choice in most of the DL robustness literature. We use CERTIFY algorithm in [4] to obtain certified radius and choose $N_0 = 100, N = 100,000, \alpha = 0.001$ in CERTIFY.

**Baseline models.** Four baseline models are discussed in this section: MACER [22], SmoothAdv [19], SmoothMix [10] and Consistency [11]. For MACER, we follow the configurations given by Table 4 in the original paper [22]. For SmoothAdv, we pick the best models under different $\sigma = 0.25, 0.50, 1.00$ from the Github repository of [19]. See Tab. 2 for more details on hyper-parameter selection of SmoothAdv. For SmoothMix and Consistency, we follow the same setting as in the original papers.

**AdvMacer .** We apply Algorithm 1 to train our **Adv-Macer** models. On Cifar-10, we choose $\gamma = 8.0, \lambda = 12.0, \beta = 16.0$ for all $\sigma = 0.25, 0.50, 1.00$. The choices of $T, m, \epsilon$ are the same as SmoothAdv to ensure fair comparison and are summarized in Tab. 2. We follow the same training scheme as [19]. The initial learning rate is 0.1 and

decays by a factor of 0.1 every 50 epochs. A batch size of 256 is used in the training. For more details, please refer to [19]. Note that by the choice of hyper-parameters, SmoothAdv and **AdvMacer** have the same training time, which implies the improved performance of **AdvMacer** is not gained from more expensive computation. The experiment results on Cifar-10 are summarized in Tab. 1.

Besides, the training time of three major models (SmoothAdv, MACER, **AdvMacer** ) is reported in Tab. 2. Compared with SmoothAdv, **AdvMacer** has the same training cost but achieves both higher accuracy and higher ACR in all $\sigma$ (as in Tab. 1). Although **AdvMacer** has more training time for each epoch than MACER does, **AdvMacer** is still more computationally efficient due to fewer required training epochs. For example, with $\sigma = 1$, MACER requires 440 epochs to obtain ACR 0.518 in 32.8 hrs, while **AdvMacer** only needs 150 epochs to achieve higher ACR 0.554 in 8 hrs. Thus, **AdvMacer** improves 7% ACR while being $4\times$ faster to train. Training efficiency of **AdvMacer** compared with MACER can be observed in Figure 1.

**EsbRs.** We also employ our proposed ensemble techniques introduced in Sec. 3.2 to enhance robustness performance. **(I).** For mixed-model ensemble, we use the following naming convention to report our result: **EsbRs-**Model1$\times$n+Model2$\times$m represents the ensemble model obtained by $n$ Model1 and $m$ Model2. For example, **EsbRs-AdvMacer** $\times$1+SmoothAdv$\times$2 represents the ensemble of one **AdvMacer** model and two SmoothAdv models. Our empirical experiments in Figure 3 also verifies the theoretical analysis on the success of mixed-model ensemble in Sec. 3.2. In ensemble experiments, we independently train all **Adv-Macer** and SmoothAdv models on Cifar-10 with $\sigma = 0.50$ and the model configuration is given by Tab. 2. In Figure 3, we observe that mixed-model ensemble gives universally better ACR than intra-model ensemble (both perform better than non-ensemble baseline), which is in accordance to the analysis in Sec. 3.2. **(II).** It can be seen from Tab. 3 that all ensemble strategies increase ACR. Max margin ensemble (MME) [21] slightly underperforms average weighted **EsbRs**, while our optimal weighted **EsbRs** can increase ACR by up to 4%.

**Experiment on ImageNet, SVHN and CIFAR** We report accuracy and ACR of SmoothAdv, MACER and **AdvMacer** on ImageNet with $\sigma = 0.25/0.50/1.00$. We train MACER model using the configuration provided by Table 5 of [22]. For SmoothAdv, we choose the best configuration from [19], that is $T = 1, m = 1, \epsilon = 0.5$ and 90 epochs. For **AdvMacer** , we use the same configuration as SmoothAdv and set $\lambda = 2, \gamma = 8, \beta = 16$ for these additional parameters. The performance and training cost is reported in Tab. 4, from which we see that **AdvMacer** achieves best ACR in all $\sigma$ while being the most cost-efficient. The additional experiments results on SVHN are reported in Appendix D. The cer-

Table 1: Cifar-10: ACR on 500 test images of Cifar-10. Clean accuracy is reported in parenthesis.

| | Methods | $\sigma = 0.25$ | $\sigma = 0.5$ | $\sigma = 1.0$ | Ensemble? |
|---|---|---|---|---|---|
| Baselines | SmoothAdv [19] | 0.541 (74.2%) | 0.735 (56.4%) | 0.758 (45.8%) | × |
| | MACER [22] | 0.518 (79.4%) | 0.682 (63.4%) | 0.768 (42.4%) | × |
| | SmoothMix [10] | 0.545 (76.0%) | 0.685 (63.8%) | 0.626 (48.4%) | × |
| | Consistency [11] | 0.535 (78.4%) | 0.701 (64.6%) | 0.719 (45.8%) | × |
| Ours | **AdvMacer** | **0.554** (76.0%) | **0.742** (58.4%) | **0.794** (47.6%) | × |
| | **EsbRs-AdvMacer** ×3 | **0.583** (76.4%) | 0.772 (58.8%) | 0.805 (47.6%) | √ |
| | **EsbRs-AdvMacer** ×1+SmoothAdv×2 | 0.572 (77.2%) | **0.783** (59.4%) | **0.810** (47.2%) | √ |
| | **EsbRs-AdvMacer** ×2+MACER×1 | 0.568 (79.8%) | 0.728 (63.6%) | 0.801 (42.8%) | √ |
| | **EsbRs-AdvMacer** ×1+MACER×2 | 0.570 (80.4%) | 0.723 (65.0%) | 0.760 (44.0%) | √ |

Table 2: Model configuration and training cost: main hyper-parameters and training time for SmoothAdv, MACER and **AdvMacer** on Cifar-10 with varing $\sigma$. For the additional parameters in **AdvMacer**, we pick $\lambda = 12.0, \gamma = 8.0, \beta = 16.0$. ACR is also reported for easy comparison.

| Models | $T$ | $m$ | $\epsilon$ | Epochs | ACR | Time |
|---|---|---|---|---|---|---|
| SmoothAdv/**AdvMacer** ($\sigma = 0.25$) | 2 | 8 | 1.0 | 150 | 0.541/**0.554** | 15.5h |
| SmoothAdv/**AdvMacer** ($\sigma = 0.50$) | 2 | 8 | 2.0 | 150 | 0.735/**0.742** | 15.5h |
| SmoothAdv/**AdvMacer** ($\sigma = 1.00$) | 2 | 4 | 2.0 | 150 | 0.758/**0.794** | 8h |
| MACER ($\sigma = 0.25/0.50/1.00$) | NA | 16 | NA | 440 | 0.518/0.682/0.768 | 32.8h |

Table 3: Optimal weighted 2-ensemble. ACR and clean accuracy on 500 test images of Cifar-10 with $\sigma = 0.25$. All certification has parameters $N_0 = 100, N = 100,000, \alpha = 0.001$. Optimal weights are computed from Algorithm 2 with **AdvMacer** models $F^1, F^2$, $m = 10, n = 10, t = 0.3, \sigma = 0.25, \tilde{\sigma} = 0.01$.

| Model | Accuracy | ACR | Certificate Time |
|---|---|---|---|
| **AdvMacer** (no ensemble) | 0.760 | 0.554 | 8.9s |
| Avg wt (baseline) | 0.760 | 0.572 | 18.0s |
| MME [21] (baseline) | 0.754 | 0.567 | 19.6s |
| **Optimal weight** (ours) | **0.766** | **0.576** | 26.3s |

tified accuracy table for Cifar-10 under $\sigma = 0.25/0.50/1.00$ can be found in Appendix C.

Table 4: ImageNet: ACR on 500 test images of ImageNet. Clean accuracy is reported in parenthesis.

| Methods | $\sigma = 0.25$ | $\sigma = 0.5$ | $\sigma = 1.0$ | Time |
|---|---|---|---|---|
| SmoothAdv | 0.519 (61.5%) | 0.801 (55.6%) | 0.971 (41.4%) | 48h |
| MACER | 0.438 (63.2%) | 0.628 (52.6%) | 0.634 (37.8%) | 70h |
| **AdvMacer** (ours) | **0.537** (63.9%) | **0.837** (56.2%) | **0.989** (45.6%) | 48h |

**Performance and Discussion. (I) Better performance.** From Tab. 1, Tab. 2, Tab. 4 and Tab. 8 in Appendix D,

we can conclude that our **AdvMacer** gives the best ACR among all $\sigma$'s and all SOTA baselines (SmoothAdv, MACER, SmoothMix, Consistency) on various datasets (Cifar-10, ImageNet, SVHN). **(II) Higher efficiency.** In terms of training cost, **AdvMacer** is 4x faster than MACER (while still achieving higher ACR); as fast as the most efficient SmoothAdv models but achieves higher ACR; 2.5x faster than Consistency (while still achieving larger ACR). **(III) Mixed-model ensemble.** For ensemble models, **AdvMacer** ×1+SmoothAdv×2 outperforms all the other models on Cifar-10 with $\sigma = 0.50, 1.00$, suggesting that one may prefer mixed-model ensemble in particular situations. Different from [9], we are the first to observe mixed-model ensemble can outperform intra-model ensemble and perform analysis to explain this opposite phenomenon. Besides, the introduction of **AdvMacer** brings enriched diversity of component models, greatly improving the power of mixed-model ensemble and making **AdvMacer** and **EsbRs** *integrated* contributions.

## 5. Conclusions

We have proposed two novel and cost-effective approaches to promote robustness of randomized smoothed classifiers: **AdvMacer** improve the robustness by maximizing the certified radius over adversarial example, and **EsbRs** can further improve **AdvMacer** on both clean accuracy and robustness certificate. Through extensive numerical experiments, we show that **AdvMacer** outperforms major baseline

models (SmoothAdv, MACER, Consistency, SmoothMix) on various datasets (Cifar-10, ImageNet, SVHN). The ACR improvement is up to 15% compared with MACER and 8% compared with the best models of SmoothAdv. Moreover, we provided a general theoretical analysis for **EsbRs** and develop a theoretical-grounded methodology to design optimal ensemble scheme, which outperforms prior works.

## Acknowledgement

## References

[1] Motasem Alfarra, Adel Bibi, Philip HS Torr, and Bernard Ghanem. Data dependent randomized smoothing. In *Uncertainty in Artificial Intelligence*, pages 64–74. PMLR, 2022. 4

[2] Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020. 3

[3] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pages 274–283. PMLR, 2018. 1

[4] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019. 1, 2, 5, 7, 13

[5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 7

[6] Francisco Eiras, Motasem Alfarra, M Pawan Kumar, Philip HS Torr, Puneet K Dokania, Bernard Ghanem, and Adel Bibi. Ancer: Anisotropic certification via sample-wise volume maximization. *arXiv preprint arXiv:2107.04570*, 2021. 4

[7] I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015. 1

[8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 7

[9] Miklós Z Horváth, Mark Niklas Mueller, Marc Fischer, and Martin Vechev. Boosting randomized smoothing with variance reduced classifiers. In *International Conference on Learning Representations*. 3, 4, 5, 6, 7, 8

[10] Jongheon Jeong, Sejun Park, Minkyu Kim, Heung-Chang Lee, Do-Guk Kim, and Jinwoo Shin. Smoothmix: Training confidence-calibrated smoothed classifiers for certified robustness. *Advances in Neural Information Processing Systems*, 34, 2021. 3, 7, 8

[11] Jongheon Jeong and Jinwoo Shin. Consistency regularization for certified robustness of smoothed classifiers. *Advances in Neural Information Processing Systems*, 33:10558–10570, 2020. 3, 7, 8

[12] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 7

[13] M. Lecuyer, V. Atlidakis, R. Geambasu, D. Hsu, and S. Jana. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 656–672, 2019. 1

[14] Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive noise. In *NeurIPS*, 2019. 1

[15] Y. Li, X. Bian, and S. Lyu. Attacking object detectors via imperceptible patches on background. *arXiv preprint arXiv:1809.05966*, 2018. 1

[16] Chizhou Liu, Yunzhen Feng, Ranran Wang, and Bin Dong. Enhancing certified robustness via smoothed weighted ensembling. *arXiv preprint arXiv:2005.09363*, 2020. 4

[17] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. 1, 2

[18] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 7

[19] Hadi Salman, Greg Yang, Jerry Li, Pengchuan Zhang, Huan Zhang, Ilya Razenshteyn, and Sébastien Bubeck. Provably robust deep learning via adversarially trained smoothed classifiers. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 11292–11303, 2019. 1, 2, 4, 7, 8, 13, 14

[20] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *ICLR*, 2014. 1

[21] Zhuolin Yang, Linyi Li, Xiaojun Xu, Bhavya Kailkhura, Tao Xie, and Bo Li. On the certified robustness for ensemble models and beyond. In *International Conference on Learning Representations*, 2022. 3, 4, 5, 7, 8

[22] 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. In *International Conference on Learning Representations*, 2019. 1, 2, 3, 4, 7, 8, 13

<div align="center">**Appendix**</div>

# A. Full algorithm of AdvMacer

---

**Algorithm 1** Our **AdvMacer** $(\sigma, m, T, \lambda, \beta, \gamma)$

---

**Input:** training set $\hat{p}_{\text{data}}$, noise level $\sigma$, number of Gaussian samples $m$, regularization parameter $\lambda$, hinge factor $\gamma$, inverse temperature $\beta$, number of PGD step $T$

**for** each iteration **do**

    1) Sample a mini-batch $(x_1, y_1), \ldots, (x_n, y_n) \sim \hat{p}_{\text{data}}$

    2) For each $(x_i, y_i)$, use $T$-step SmoothAdv to generate adversarial example $\hat{x}_i$

    3) For each $(\hat{x}_i, y_i)$, draw $m$ i.i.d. Gaussian samples $x_{i1}, \ldots, x_{im}$ from $\mathcal{N}(x_i, \sigma^2 I)$

    4) Obtain an estimation of $G_\theta(\hat{x})$ by $\hat{z}_\theta(\hat{x}) \leftarrow \frac{1}{m} \sum_{k=1}^{m} F_\theta(\hat{x}_{ik})$, for $i = 1, \ldots, n$

    5) Collect the set of data with correct prediction: $\mathcal{S}_\theta = \{i : y_i = \arg\max_c \hat{z}_\theta(\hat{x}_i)_c\}$

    6) For each $i \in \mathcal{S}_\theta$, compute the second most likely class $\hat{y}_i \leftarrow \arg\max_{c \neq y_i} \hat{z}_\theta(\hat{x}_i)_c$

    7) For each $i \in \mathcal{S}_\theta$, compute $\hat{\xi}(\hat{x}_i, y_i) \leftarrow \Phi^{-1}(\hat{z}_\theta(x)_{y_i}) - \Phi^{-1}(\hat{z}_\theta(x)_{\hat{y}_i})$

    8) Sample $\delta \sim \mathcal{N}(0, \sigma^2 I)$ and update $\theta$ with SGD to minimize

$$-\frac{1}{n} \sum_{i=1}^{n} \log \hat{z}_\theta(\hat{x}_i + \delta)_{y_i} + \frac{\lambda\sigma}{2n} \sum_{i \in \mathcal{S}_\theta} \max\{\gamma - \hat{\xi}_\theta(\hat{x}_i + \delta, y_i), 0\}$$

**end for**

**Output:** model parameters $\theta$

---

# B. Full algorithm of optimal weight design

---

**Algorithm 2** Optimal weights of ensemble with 2 models. The function **ComputeWeight** will return the optimal weights

---

**Input:** two base models $F^1, F^2$, number of Gaussian noise $n$, number of perturbed models $m$, noise level $\sigma$, proportion of the parameters to perturb $t$, standard deviation of perturbation on parameters $\tilde{\sigma}$, query point $x$, target label $y$

**function** PerturbModel($F^1, F^2, m, t, \tilde{\sigma}$):

    **for** $l = 1, 2$ **do**

        **for** each $j = 1, \ldots, m$ **do**

            **for** each parameter $\theta$ in $F^l$ **do**

                Draw a Bernoulli variable $X$ from Bernoulli($t$)

                **if** $X = 1$ **then**

                    Draw $\delta \sim \mathcal{N}(0, \tilde{\sigma}^2)$ and update $\theta \leftarrow \theta + \delta$

                **end if**

            **end for**

            Store perturbed model $\hat{F}_j^l$

        **end for**

    **end for**

    **Output:** perturbed models $\hat{F}_1^1, \ldots, F_m^1$ and $\hat{F}_1^2, \ldots, \hat{F}_m^2$.

**function** Estimation($[\hat{F}_1^1, \ldots, F_m^1], [\hat{F}_1^2, \ldots, \hat{F}_m^2], \sigma, n, x, y$):

    Draw $n$ i.i.d. noisy samples from $\mathcal{N}(x, \sigma^2 I)$ and denote them by $x_1, \ldots, x_n$

    **for** $l = 1, 2$ **do**

        **for** $i = 1, \ldots, m$ **do**

            **for** $j = 1, \ldots, n$ **do**

                Compute $z^{l,(i-1)n+j} \leftarrow \hat{F}_i^l(x_j)_y \mathbf{1} - \hat{F}_i^l(x_j)$, where $\mathbf{1} = [1, 1, \ldots, 1]^\top$ is the vector of all 1's.

            **end for**

        **end for**

    **end for**

    **Output:** estimates of logits $z^{l,j}$ for $l = 1, 2$ and $j = 1, \ldots, mn$

**function** ComputeWeight($F^1, F^2, \sigma, \tilde{\sigma}, n, m, t, x, y$):

    1) $[\hat{F}_1^1, \ldots, F_m^1], [\hat{F}_1^2, \ldots, \hat{F}_m^2] \leftarrow$ PerturbModel($F^1, F^2, m, t, \tilde{\sigma}$)

    2) $z^{l,j} \leftarrow$ Estimation($[\hat{F}_1^1, \ldots, F_m^1], [\hat{F}_1^2, \ldots, \hat{F}_m^2], \sigma, n, x, y$) for $l = 1, 2$ and $j = 1, \ldots, m$

    3) Compute $\bar{z}^l \leftarrow \frac{1}{mn} \sum_{j=1}^{mn} z^{l,j}$ for $l = 1, 2$ and $a_i \leftarrow \min\{\bar{z}_i^1, \bar{z}_i^2\}^{-2}$ for $i = 1, \ldots, c$

    4) Compute

$$b_i \leftarrow \frac{1}{mn} \sum_{j=1}^{mn} (z^{1,j} - \bar{z}^1)(z^{1,j} - \bar{z}^1)_{ii}^\top, \ c_i \leftarrow \frac{2}{mn} \sum_{j=1}^{mn} (z^{1,j} - \bar{z}^1)(z^{2,j} - \bar{z}^2)_{ii}^\top, \ d_i \leftarrow \frac{1}{mn} \sum_{j=1}^{mn} (z^{2,j} - \bar{z}^2)(z^{2,j} - \bar{z}^2)_{ii}^\top,$$

    5) Compute

$$A = \sum_{i=2}^{c} a_i(b_i + c_i + d_i), \quad B = \sum_{i=2}^{c} -a_i(c_i + 2d_i), \quad C = \sum_{i=2}^{c} a_i d_i.$$

    **if** $A > 0$ and $0 \leq -\frac{B}{2A} \leq 1$ **then**

        $w_1 \leftarrow -\frac{B}{2A}$ and $w_2 \leftarrow 1 + \frac{B}{2A}$

    **else**

        $w_1 \leftarrow 0, w_2 \leftarrow 1$ if $A + B > 0$ else $w_1 \leftarrow 1, w_2 \leftarrow 0$

    **end if**

    **Output:** $w_1, w_2$

---

## C. Certified accuracy

We also provide certified accuracy table for Cifar-10, which is presented in Tab. 5, Tab. 6 and Tab. 7.

Table 5: Certified accuracy: certified accuracy and ACR of the first 500 test images of Cifar-10 with $\sigma = 0.25$. Each column represents the robust accuracy that can be certified at this $\ell_2$ radius.

| Model ($\sigma = 0.25$) | 0.00 | 0.25 | 0.50 | 0.75 | 1.00 | 1.25 | 1.50 | 1.75 | 2.00 | ACR |
|---|---|---|---|---|---|---|---|---|---|---|
| SmoothAdv | 0.742 | 0.660 | **0.572** | 0.45 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.541 |
| MACER | **0.794** | **0.678** | 0.524 | 0.400 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.518 |
| SmoothMix | 0.76 | 0.688 | 0.572 | 0.446 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.545 |
| **AdvMacer** | 0.76 | 0.668 | **0.572** | **0.484** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **0.554** |
| **EsbRs-AdvMacer** ×3 | 0.764 | 0.700 | **0.614** | 0.514 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **0.583** |
| **EsbRs-**SmoothAdv×3 | 0.766 | 0.698 | 0.600 | 0.506 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.576 |
| **EsbRs-AdvMacer** ×1+SmoothAdv×2 | 0.772 | 0.672 | 0.594 | 0.498 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.572 |
| **EsbRs-AdvMacer** ×2+MACER×1 | 0.798 | 0.700 | 0.586 | 0.472 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.568 |
| **EsbRs-AdvMacer** ×1+MACER×2 | **0.804** | **0.714** | 0.598 | 0.462 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.570 |

Table 6: Certified accuracy: certified accuracy and ACR of the first 500 test images of Cifar-10 with $\sigma = 0.50$. Each column represents the robust accuracy that can be certified at this $\ell_2$ radius.

| Model ($\sigma = 0.50$) | 0.00 | 0.25 | 0.50 | 0.75 | 1.00 | 1.25 | 1.50 | 1.75 | 2.00 | ACR |
|---|---|---|---|---|---|---|---|---|---|---|
| SmoothAdv | 0.564 | 0.516 | 0.468 | 0.432 | 0.394 | 0.328 | **0.286** | **0.224** | 0.0 | 0.735 |
| MACER | **0.634** | **0.566** | 0.476 | 0.432 | 0.346 | 0.258 | 0.206 | 0.126 | 0.0 | 0.682 |
| SmoothMix | 0.638 | 0.548 | 0.48 | 0.416 | 0.34 | 0.274 | 0.21 | 0.152 | 0.0 | 0.685 |
| **AdvMacer** | 0.584 | 0.532 | **0.486** | **0.442** | **0.398** | **0.334** | 0.270 | 0.216 | 0.0 | **0.742** |
| **EsbRs-AdvMacer** ×3 | 0.588 | 0.544 | 0.498 | 0.448 | 0.414 | 0.360 | 0.288 | 0.230 | 0.0 | 0.772 |
| **EsbRs-**SmoothAdv×3 | 0.584 | 0.530 | 0.476 | **0.454** | 0.420 | 0.362 | 0.308 | **0.254** | 0.0 | 0.777 |
| **EsbRs-AdvMacer** ×1+SmoothAdv×2 | 0.594 | 0.540 | 0.482 | **0.454** | **0.422** | **0.374** | **0.310** | 0.238 | 0.0 | **0.783** |
| **EsbRs-AdvMacer** ×2+MACER×1 | 0.636 | 0.564 | 0.504 | 0.446 | 0.388 | 0.294 | 0.240 | 0.172 | 0.0 | 0.728 |
| **EsbRs-AdvMacer** ×1+MACER×2 | **0.650** | **0.568** | **0.506** | 0.450 | 0.370 | 0.292 | 0.220 | 0.158 | 0.0 | 0.723 |

Table 7: Certified accuracy: certified accuracy and ACR of the first 500 test images of Cifar-10 with $\sigma = 1.00$. Each column represents the robust accuracy that can be certified at this $\ell_2$ radius.

| Model ($\sigma = 1.00$) | 0.00 | 0.25 | 0.50 | 0.75 | 1.00 | 1.25 | 1.50 | 1.75 | 2.00 | 2.25 | ACR |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SmoothAdv | 0.458 | 0.418 | 0.374 | 0.312 | 0.288 | 0.254 | 0.234 | 0.196 | 0.180 | 0.158 | 0.758 |
| MACER | 0.424 | 0.392 | 0.354 | 0.328 | **0.304** | **0.274** | **0.250** | **0.212** | 0.184 | 0.156 | 0.768 |
| SmoothMix | 0.484 | 0.42 | 0.348 | 0.292 | 0.244 | 0.22 | 0.178 | 0.15 | 0.124 | 0.096 | 0.626 |
| **AdvMacer** | **0.476** | **0.440** | **0.392** | **0.350** | 0.302 | **0.274** | 0.236 | **0.212** | **0.186** | **0.164** | **0.794** |
| **EsbRs-AdvMacer** ×3 | **0.476** | 0.426 | 0.386 | 0.358 | 0.302 | 0.270 | 0.246 | 0.220 | 0.196 | 0.172 | 0.805 |
| **EsbRs-**SmoothAdv×3 | 0.466 | **0.432** | 0.388 | **0.360** | 0.294 | 0.260 | 0.240 | 0.212 | 0.186 | 0.170 | 0.801 |
| **EsbRs-AdvMacer** ×1+SmoothAdv×2 | 0.472 | **0.432** | **0.394** | 0.356 | 0.304 | 0.262 | 0.240 | 0.214 | 0.194 | **0.174** | **0.810** |
| **EsbRs-AdvMacer** ×2+MACER×1 | 0.428 | 0.404 | 0.370 | 0.342 | **0.318** | 0.276 | **0.256** | **0.230** | **0.198** | 0.164 | 0.801 |
| **EsbRs-AdvMacer** ×1+MACER×2 | 0.440 | 0.404 | 0.366 | 0.340 | 0.304 | **0.282** | 0.238 | 0.202 | 0.184 | 0.154 | 0.760 |

## D. Additional experiments on SVHN dataset

We compare the performance of SmoothAdv, MACER and **AdvMacer** on SVHN dataset with $\sigma = 0.25, 0.50$. On SVHN with $\sigma = 0.25$, we choose $T = 2$, $m = 4$, $\lambda = 12.0$, $\gamma = 8.0$, $\beta = 16.0$, $\epsilon = 0.5$ and train the model for 150 epochs. On SVHN with $\sigma = 0.50$, we still choose $T = 2$, $m = 4$, $\gamma = 8.0$, $\beta = 16.0$, $\epsilon = 0.5$ but a different $\lambda = 4.0$. The model is also trained for 150 epochs. The initial learning rate is set to 0.01 and drops by a factor of 0.1 every 50 epochs. The other training details follow the same as Cifar-10. For SmoothAdv, take $T = 2$, $m = 4$, $\epsilon = 0.5$ when $\sigma = 0.25$ and $T = 2$, $m = 4$, $\epsilon = 0.25$ when $\sigma = 0.50$. We train MACER model for 440 epochs whose configuration is given by C.2.2 of [22]. We report the experiment results in Tab. 8.

Table 8: SVHN: clean accuracy and ACR of different models evaluated on 500 test images of SVHN with varing $\sigma$.

| $\sigma$ | Model | Accuracy | ACR | Training Time |
|---|---|---|---|---|
| | SmoothAdv | 85.8% | 0.560 | 11.4h |
| | MACER | 86.8% | 0.549 | 48.5h |
| | **AdvMacer** | 86.6% | **0.569** | 11.4h |
| 0.25 | **EsbRs-**SmoothAdv×3 | 87.8% | 0.578 | NA |
| | **EsbRs-AdvMacer** ×3 | 88.2% | **0.582** | NA |
| | **EsbRs-AdvMacer** ×1+MACER×2 | 87.8% | 0.559 | NA |
| | **EsbRs-AdvMacer** ×2+MACER×1 | 88.6% | 0.570 | NA |
| | **EsbRs-AdvMacer** ×1+SmoothAdv×2 | 87.8% | 0.577 | NA |
| | **EsbRs-AdvMacer** ×2+SmoothAdv×1 | 87.6% | **0.582** | NA |
| | SmoothAdv | 71.2% | 0.552 | 11.4h |
| | MACER | 58.4% | 0.535 | 48.5h |
| | **AdvMacer** | 67.8% | **0.572** | 11.4h |
| 0.50 | **EsbRs-**SmoothAdv×3 | 71.2% | 0.573 | NA |
| | **EsbRs-AdvMacer** ×3 | 70.4% | **0.588** | NA |
| | **EsbRs-AdvMacer** ×1+MACER×2 | 62.8% | 0.551 | NA |
| | **EsbRs-AdvMacer** ×2+MACER×1 | 66.0% | 0.564 | NA |
| | **EsbRs-AdvMacer** ×1+SmoothAdv×2 | 71.8% | 0.577 | NA |
| | **EsbRs-AdvMacer** ×2+SmoothAdv×1 | 71.2% | 0.583 | NA |

## E. Proof of Theorem 3.1

For completeness, recall that we have $k$ trained soft classifiers $F^1, \ldots, F^k : \mathbb{R}^d \to P(\mathcal{Y})$ and $\mathcal{Y} = \{1, \ldots, c\}$. Consider soft-ensemble model $H$ whose output is a weighted average of the probabilities from $F^1, \ldots, F^k$: $H(x) = \sum_{l=1}^{k} w_l(x) F^l(x)$. Suppose the associated hard classifier is $h(x) = \arg\max_{c \in \mathcal{Y}} \left( H(x) \right)_c$. Then we apply RS to $h$ and get the corresponding smoothed classifier $e(x) = \arg\max_c \mathbb{P}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)}(h(x + \epsilon) = c)$. It can be shown that the RS classifier $e$ of an ensemble model also has the same certification guarantee as in [4], which is presented in the following theorem.

**Theorem E.1** (Robustness gaurantee for ensemble). *Suppose that under Gaussian perturbation $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$, the most likely class $c_A$ is returned by the **EsbRs** $e$ with probability $p_A$ and the second most likely class $c_B$ is returned with probability $p_B$, then we have $e(x + \delta) = e(x)$ for all $\|\delta\|_2 \leq R$, where $R = \frac{\sigma}{2}(\Phi^{-1}(p_A) - \Phi^{-1}(p_B))$.*

Without loss of generality, we assume $\sigma = 1$. The proof follows similarly for $\sigma \neq 1$. Before presenting the formal proof, we proceed with several preparatory lemmas from [19].

**Theorem E.2** (Lemma 2 from [19]). *Let $f : \mathbb{R}^n \to [0, 1]$ be a deterministic (non-random) function and define $\hat{f}$ by*

$$\hat{f}(x) = (f * \mathcal{N}(0, I))(x) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} f(t) \exp\left( -\frac{1}{2}\|x - t\|^2 \right) dt.$$

*Let $\Phi(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{a} \exp\left( -\frac{1}{2}s^2 \right) ds$. Then the map $x \mapsto \Phi^{-1}(\hat{f}(x))$ is 1-Lipschitz.*

We refer the readers to [19] for the proof of Theorem E.2. Now we are ready to prove Theorem 3.1.

*Proof of Theorem 3.1.* Let $f_i(x) = \mathbb{I}(e(x) = i)$ for $i \in \mathcal{Y} = \{1, 2, \ldots, c\}$. Define

$$\hat{f}_i(x) = (f_i * \mathcal{N}(0, I))(x) = \mathbb{P}_{\epsilon \sim \mathcal{N}(0,I)}(e(x + \delta) = c_i).$$

From this definition, we immediately have $p_A = \hat{f}_A(x)$ and $p_B = \hat{f}_B(x)$. Then by Theorem E.2, we have $\Phi^{-1}(\hat{f}_i(x))$ is 1-Lipschitz for all $i \in \{1, \ldots, c\}$. So for any $\delta$,

$$\Phi^{-1}(\hat{f}_A(x)) - \Phi^{-1}(\hat{f}_A(x + \delta)) \leq \|\delta\|$$

Suppose $\delta^*$ is a successful adversarial noise, namely $\hat{f}_A(x + \delta^*) \leq \hat{f}_B(x + \delta^*)$ for some class $c_B$, then

$$\Phi^{-1}(\hat{f}_A(x)) - \Phi^{-1}(\hat{f}_B(x + \delta^*)) \leq \Phi^{-1}(\hat{f}_A(x)) - \Phi^{-1}(\hat{f}_A(x + \delta^*)) \leq \|\delta^*\|. \tag{11}$$

Apply Theorem E.2 again to $\hat{f}_B$ and we obtain

$$\Phi^{-1}(\hat{f}_B(x + \delta^*)) - \Phi^{-1}(\hat{f}_B(x)) \leq \|\delta^*\|. \tag{12}$$

Combining Equation (11) and Equation (12), we have that if $\delta^*$ is a successful attack, it has to satisfy

$$\|\delta^*\| \geq \frac{1}{2}(\Phi^{-1}(p_A) - \Phi^{-1}(p_B)).$$

In other words, we are guaranteed that $e(x) = e(x + \delta)$ if

$$\|\delta\|_2 \leq \frac{1}{2}(\Phi^{-1}(p_A) - \Phi^{-1}(p_B)).$$

$\square$