

CERTIFIED DEFENSE AGAINST COMPLEX ADVERSARIAL ATTACKS WITH DYNAMIC SMOOTHING

Anonymous authors

Paper under double-blind review

ABSTRACT

Randomized smoothing has emerged as a certified defence mechanism with probabilistic guarantees that works at scale. However, current randomized smoothing methods offer theoretical guarantees that are limited by their reliance on specific noise distributions, and they struggle to handle complex adversarial attacks. In this paper, we propose a novel certification method based on randomized smoothing designed to handle complex adversarial attacks, including combinations of multiple attack types. We call this method Dynamic Smoothing (DSMOOTH). Our key idea is to incorporate more general distributions for smoothing than isotropic Gaussian noise, for which probabilistic guarantees can be derived in terms of the Mahalanobis distance. These general distributions make the smoothed classifier more robust against a wide range of threats, including localized adversarial attacks and multi-attacks. We validate the performance of our method experimentally on challenging threat models using CIFAR-10 and IMAGENET, and demonstrate its superiority over state-of-the-art defenses in terms of certified accuracy. Our results show that the proposed method significantly improves the robustness of machine learning models against complex attacks, advancing their suitability for use in safety-critical applications. Code: [removed for review]

1 INTRODUCTION

Machine Learning has seen considerable progress in recent years, especially with deep neural networks (DNNs). However, these networks are vulnerable to adversarial examples (Szegedy et al., 2013; Goodfellow et al., 2014; Zhao et al., 2023), posing a challenge for their use in safety-critical areas (Kurakin et al., 2016; Shayegani et al., 2023). Adversarial attacks, such as DeepFool (Moosavi-Dezfooli et al., 2016), AutoAttack (Croce & Hein, 2020), patch-based attacks (Brown et al., 2017b), and attacks on LLMs (Zou et al., 2023) continue to evolve, outpacing existing defenses and creating a persistent struggle between attackers and defenders (Carlini & Wagner, 2017a; Madry et al., 2017a). Current defenses, e.g., denoising generative models (Gu & Rigazio, 2014; Ho et al., 2020), adversarial training (Miller et al., 2020; Kireev et al., 2022), and defensive distillation (Papernot et al., 2016a; Wang et al., 2021), have not fully succeeded in preventing stronger attacks. Hence, the problem of building trustworthy ML systems suitable for critical applications remains an open question.

Certified robustness has emerged as an alternative approach, with randomized smoothing (Lecuyer et al., 2019; Li et al., 2019; Cohen et al., 2019b; Anderson & Sojoudi, 2022; Scholten et al., 2023; Anani et al., 2024) being a notable method. This technique, which provides *probabilistic* guarantees, involves creating a smoothed classifier by applying Gaussian noise to the base classifier. This method was shown by Lecuyer et al. (2019) and Li et al. (2019) to provide consistent classification within a certified radius under ℓ_2 norm considerations, although the guarantees were initially loose. Cohen et al. (2019b) were the first to offer tight robustness guarantees for this method against ℓ_2 norm-constrained adversarial attacks, sparking further studies in this area.

Randomized smoothing has become a widely recognized method for certified robustness, though it has limitations. Cohen et al. (2019b) identified the need for further exploration of ℓ_p norms beyond ℓ_2 . Recent works have been addressing robustness guarantees for randomized smoothing against various types of adversaries, including ℓ_1 -bounded attacks (Teng et al., 2020), ℓ_0 -bounded attacks (Levine & Feizi, 2020c; Lee et al., 2019), and Wasserstein attacks (Levine & Feizi, 2020a). However, defending against complex, high-dimensional adversarial attacks remains an open challenge.

Our Contribution.

- We provide a certification method based on randomized smoothing, which we refer to as **Dynamic SMOOTHing** (DSMOOTH, Sec. 4.1). DSMOOTH uses more complex smoothing distributions than traditional randomized smoothing, making the smoothing process more adaptable to localized and non-uniform adversarial attacks than previous methods. Our method is also a suitable certified defense method against attacks based on multiple norms, such as multi-attacks.
- We derive probabilistic guarantees based on the Mahalanobis distance (Sec. 4.2). Our analysis, which is non-trivial, provides a framework to derive guarantees using push-forward measures (Thm. 4.4), which can be of independent interest. Furthermore, we derive probabilistic guarantees using the ℓ_2 norm, recovering known guarantees for isotopic Gaussian noise (Cor. 4.7).
- We provide extensive experiments on CIFAR-10 and IMAGENET, considering a multi-attack that combines the Square Attack algorithm (Andriushchenko et al., 2020) and FGSM (Goodfellow et al., 2015). We show that DSMOOTH achieves good certified accuracy, significantly outperforming baselines (Sec. 5).

2 RELATED WORK

Since there is a large amount of scientific articles on this topic, we only discuss the contributions relevant for this work. The interested reader can refer to, e.g., Kumari et al. (2023); Kwiatkowska & Zhang (2023), for a more complete overview. Defenses against adversarial examples fall into empirical and certified categories. Empirical defenses, e.g., adversarial training (Madry et al., 2017a;b; Jin et al., 2023), aim to enhance robustness but lack guarantees of being unbreakable, as many have been compromised by stronger attacks, e.g., (Carlini & Wagner, 2017b; Athalye et al., 2018; Tramèr et al., 2020). Certified defenses and verification methods ensure consistent classifier output within a small neighborhood of x , using exact methods, e.g., (Huang et al., 2017; Katz et al., 2017; Ehlers, 2017; Mao et al., 2023; 2024), or conservative methods, e.g., (Wong & Kolter, 2018; Raghunathan et al., 2018; Dvijotham et al., 2018). Randomized smoothing has emerged as a probabilistic certified defense mechanism that works at scale.

Although the literature on randomized smoothing largely focuses on simple threat models, such as imperceptible adversarial perturbations of the input images (Szegedy et al., 2014; Goodfellow et al., 2015; Papernot et al., 2016b; Carlini & Wagner, 2017c), more complex threat models have been considered. Patch attacks, which place imperceptible modifications on images, can cause misclassifications and compromise system security. Levine & Feizi (2020b) address this with (De-) Randomized Smoothing for certifiable defense, leveraging the constraints of patch attacks over general sparse attacks. Zhang et al. (2023) introduce DRSM (De-randomized smoothed MalConv), adapting de-randomized smoothing for malware detection through executables (Raff et al., 2018). Recently, randomized smoothing has been used against image transformations (Fischer et al., 2020). Randomized smoothing has also been extended to discrete data (Bojchevski et al., 2020).

3 FRAMEWORK

3.1 PROBLEM DESCRIPTION

We are given a pre-trained classifier f . We do not make any specific assumption on the inner workings of f . For instance, f can be a large convolutional neural network, e.g., ResNet (He et al., 2016), MobileNet (Howard et al., 2017), or any other model suitable for perception tasks in autonomous vehicles. We consider a threat model for the classifier f . This threat model generates adversarial images \hat{x} by adding perturbations $\hat{\delta}$ to input images x , with the goal of fooling the classifier at inference time. In this work, we consider general *white-box* adversarial attacks, i.e., attacks in which the attacker may have full access to and knowledge of the target model’s architecture, parameters, and training data. Formally, we consider the following class of adversarial attacks:

Definition 3.1. Consider a classifier f with a loss function \mathcal{L} . For an input x with label y , a *white-box* attack for f generates an adversarial example $\hat{x} = x + \hat{\delta}$, such that

$$\hat{\delta} = \arg \max_{\delta \in \mathcal{C}(\delta)} \mathcal{L}(f(x + \delta), y). \quad (1)$$

Here, \mathcal{L} is a loss function and $\mathcal{C}(\delta)$ is a perturbation set, which is the collection of all possible perturbations δ that can be applied to an input x . A white-box multi-attack is an adversarial strategy that combines multiple white-box attacks as in equation 1 to generate a single adversarial example.

Adversarial attacks as in equation 1 encompass a wide variety of attacks, such as spatial perturbations (Engstrom et al., 2019), Wasserstein- bounded perturbations (Hu et al., 2020; Wong et al., 2019), perturbations of the image colors (Laidlaw & Feizi, 2019) or perceptual adversarial attacks (Laidlaw et al., 2021; Wong & Kolter, 2021). Adversarial attacks on traffic sign detection by (Li et al., 2021) and physical adversarial attacks (Brown et al., 2017a; Woitschek & Schneider, 2023) are also attacks as in Def. 3.1.

Multi-attacks as in Def. 3.1 combine any of these methods, to exploit a broader range of model vulnerabilities. Multi-attacks optimize perturbations under different norms, leading to more complex, non-uniform perturbations. An example of a multi-attack, which is used for experimental comparison in this work, is a combination of the Square Attack algorithm (Andriushchenko et al., 2020) with FGSM (Goodfellow et al., 2015). This attack, which we denote as SQUARE + FGSM, first applies a Square Attack to an input image, and then it applies a FGSM attack to the resulting sample.

The research question. We study the problem of providing a certified defense mechanism against adversarial attacks as in Def. 3.1. This defense mechanism ought to be suitable to handle highly-dimensional input, such as images in datasets for vision-based perception systems of robots and autonomous driving systems.

3.2 RANDOMIZED SMOOTHING

Randomized smoothing is a technique for improving the robustness of models against adversarial attacks (Lecuyer et al., 2019; Li et al., 2019; Cohen et al., 2019b). The main principle of randomized smoothing is to transform a deterministic classifier into a probabilistic one by averaging its predictions over many noisy versions of the input. This process effectively “smooths out” the decision boundary of the classifier, making it less sensitive to input perturbations. Specifically, given a classifier f , randomized smoothing is a method for constructing a new classifier g as

$$g(x) := \arg \max_y \mathbb{P}(f(x + \varepsilon) = y) \quad \text{with} \quad \varepsilon \sim \mathbb{P}(\varepsilon).$$

Here, $\mathbb{P}(\varepsilon)$ is the *smoothing distribution* and it determines how noise is added to the input x . Typically, the smoothing distribution is a Gaussian distribution of the form $\mathbb{P}(\varepsilon) = \mathcal{N}(0, \sigma^2 I)$, with I the identity matrix and σ a user-defined scalar, although other distributions have been considered (see, e.g., (Teng et al., 2020; Levine & Feizi, 2020c; Lee et al., 2019)). Randomized smoothing provides probabilistic robustness guarantees in terms of the *certified radius*. This radius specifies a region around an input x within which the smoothed classifier’s prediction is guaranteed to be robust, with a certain probability. The region specified by the certified radius is called a *safety region*. The choice of the smoothing distribution significantly affects the robustness guarantees provided by randomized smoothing. The guarantees obtained with standard Gaussian smoothing distributions, as above, specify a safety region \mathcal{S} using ℓ_p norms, e.g., $\mathcal{S} := \{\hat{x} : \|\hat{x} - x\|_p \leq R\}$ for some radius R . These types of guarantees are suitable to certify robustness against imperceptible adversarial perturbations on the input image, such as those generated by L-BFGS (Szegedy et al., 2014), FGS (Goodfellow et al., 2015), DeepFool (Moosavi-Dezfooli et al., 2016), JSMA (Papernot et al., 2016b), or CW (Carlini & Wagner, 2017c). However, due to their reliance on global noise perturbations, guarantees based on isotropic Gaussian smoothing may be unsuitable for complex attacks that use structured and localized adversarial perturbations.

4 METHODOLOGY

4.1 OVERVIEW

We extend the randomized smoothing framework by Cohen et al. (2019a) to more complex smoothing distributions. In contrast to prior work, our framework uses *anisotropic* Gaussian noise as a smoothing distribution, i.e., a Gaussian distribution in which the variances along different dimensions of the space are not equal, which allows to handle both sparse and localized adversarial perturbations. Importantly, in Sec. 4.2 we derive probabilistic guarantees for this method that generalize previous known guarantees (Cohen et al., 2019a).

To define our smoothing framework, consider general adversarial examples of the form $\hat{x} = x + \hat{\delta}$ constructed with a general white-box (multi)attack as in Def. 3.1. We can view $\hat{\delta}$ as a random variable, where the randomness is given by the choice of the corresponding natural example x . We define the covariance matrix Σ such that its entries are

$$[\Sigma]_{i,j} := \text{COV}[\hat{\delta}_i, \hat{\delta}_j] = \mathbb{E}[(\hat{\delta}_i - \mathbb{E}[(\hat{\delta}_i)])(\hat{\delta}_j - \mathbb{E}[(\hat{\delta}_j)])], \quad (2)$$

with $\hat{\delta}_i$ and $\hat{\delta}_j$ the i -th and j -th entries of the random variable $\hat{\delta}$. For an input x of dimension d , our smoothed classifier is defined as follows

$$g(x) := \arg \max_y \mathbb{P}(f(x + \delta) = y) \quad \text{with } \varepsilon \sim \mathcal{N}\left(0, \frac{\sigma^2}{\sqrt{d \det(\Sigma)}} \Sigma\right)^1 \quad (3)$$

We refer to this algorithm as **Dynamic SMOOTHing** (DSMOOTH). This algorithm dynamically adapts to adversarial attacks, since the matrix Σ embeds information on the adversarial perturbations $\hat{\delta}$. In equation 3, σ is a user-defined parameter of the smoothed classifier. As in the original work by (Cohen et al., 2019b), the parameter σ regulates the trade-off between robustness and accuracy. In fact adding more noise (a higher σ) tends to increase the robustness of the model to adversarial attacks, as the model’s predictions become more invariant to small perturbations in the input. However, this can also degrade the model’s accuracy on clean, unperturbed inputs because the predictions become more uncertain. In App. F we show examples of CIFAR-10 (Fig. 6) and IMAGENET (Fig. 7) images corrupted with the smoothing distribution as in equation 3 for a SQUARE + FGSM attack as described in Sec. 3.1. We remark that DSMOOTH as in equation 3 is essentially a generalization of the framework by Cohen et al. (2019a). In fact, by setting $\Sigma = I$ in equation 3, DSMOOTH is equivalent to the randomized smoothing algorithm in equation 1 of Cohen et al. (2019a).

Practical implementation of the smoothing algorithm as in equation 3. In general, the matrix Σ in equation 3 is unknown and it has to be learned from samples. We approximate Σ is to gather sample perturbations $\hat{\delta}$ in simulation, and then compute the resulting sample covariance matrix as in equation 2. However, the resulting smoothing algorithm as in equation 3 may be impractical when dealing with large input, since the size of Σ grows with the input size.

To overcome this problem, we use Principal Component Analysis (PCA) (Abdi & Williams, 2010) to provide a surrogate Σ_k of reduced size for the covariance matrix Σ , and sample ε as in equation 3 using Σ_k . To generate Σ_k , we use a rank- k approximation, where $k < \dim(\Sigma)$. This is done by retaining only the top k eigenvectors corresponding to the largest k eigenvalues. The approximated covariance matrix Σ_k can then be expressed as $\Sigma_k = U_k \Lambda_k U_k^T$, where U_k is a matrix of size $d \times k$ containing the top k eigenvectors, and Λ_k is a diagonal matrix of size $k \times k$ containing the top k eigenvalues. In Sec. 5 we show empirically that different choices for k do not significantly affect the performance of DSMOOTH.

Algorithms for the evaluation and certification of g as in equation 3 are given in App. A.

4.2 CERTIFICATION GUARANTEES

We derive certification guarantees for a smoothed classifier as in equation 3. In this section, we provide guarantees based on the Mahalanobis distance, which can be seen as a generalization of the ℓ_2 norm. However, we derive guarantees for our method in terms of the ℓ_2 norm in Sec. 4.3. Formally, we consider the following distance in our analysis.

Definition 4.1 (Mahalanobis Distance). *Consider adversarial examples of the form $\hat{x} = x + \hat{\delta}$ as defined in Sec. 3.1, and denote with Σ be the covariance matrix of the r.v. $\hat{\delta}$. Then, the Mahalanobis distance of \hat{x} with respect to (w.r.t.) x is defined as*

$$\text{MAHL}(\hat{x} \mid x) := \sqrt{(\hat{x} - x)^T \Sigma^{-1} (\hat{x} - x)},$$

where Σ^{-1} is the inverse of Σ .

In contrast to the standard ℓ_p norms, the Mahalanobis distance in Def. 4.1 adjusts for the spread and orientation of the adversarial perturbations δ . Since the eigenvalues of Σ are proportional to the amount of variance captured by each principal component, then any safety boundary of the

¹Throughout this work we assume that $\det(\Sigma) \neq 0$, i.e., we assume that Σ is positive-definite.

form $\text{MAHL}(\hat{x} \mid x) \leq R$ is an ellipsoidal “stretched” in the direction of the worst-case adversarial examples. In Sec. 5 we show experimentally that certified radii based on the Mahalanobis distance are better suited for complex adversarial attacks than certification bounds based on ℓ_p norm. There is a natural connection between the Mahalanobis and the ℓ_2 norm, as discussed in Sec. 4.3.

A general framework for the push-forward measure. Before discussing these results, however, we prove a general theoretical result, which is essential to provide guarantees for our proposed certification method. This result, which could be of independent interest, ensures general certification guarantees when the smoothing distribution is the push-forward distribution of an isotopic Gaussian distribution. Recall that the push-forward measure is defined as follows.

Definition 4.2 (Push-forward measure). *Given a measurable space (X, \mathcal{A}) and a measurable function $p : X \rightarrow Y$ mapping from X to Y , and a measure μ on X , the push-forward measure $p^\# \mu$ on Y is defined as $(p^\# \mu)(B) = \mu(p^{-1}(B))$ for any measurable set $B \subseteq Y$.*

In other words, sampling from the push-forward measure $p^\# \mu$ consists of first sampling from μ , and then applying the function p to this sample. In the following theorem, we extend guarantees for randomized smoothing to a generic sampling distribution of the form $p^\# \mathcal{N}(0, \sigma^2 I)$. Throughout this section, we denote with Φ^{-1} the inverse of the standard Gaussian CDF. The following theorem holds.

Theorem 4.3 (Randomized smoothing for the push-forward measure). *Consider a classifier f , and let p be a deterministic invertible function. Consider the mapping $g(x) := \arg \max_y \mathbb{P}(f(x + \delta) = y)$ with $\delta \sim p^\# \mathcal{N}(0, \sigma^2 I)$. For a class $y_A \in Y$ suppose that there exist two constants $\underline{p}_A, \overline{p}_B \in [0, 1]$ such that*

$$\mathbb{P}(f(x + \varepsilon) = y_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{y_B \neq y_A} \mathbb{P}(f(x + \varepsilon) = y_B),$$

with $\varepsilon \sim p^\# \mathcal{N}(0, \sigma^2 I)$. Then, it holds $g(x + \delta) = y_A$ for all δ such that

$$\|p^{-1}(\delta)\|_2 \leq \frac{\sigma}{2} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)).$$

Thm. 4.3 provides a defensive method that ensures guarantees in terms of the ℓ_2 norm with respect to $p^{-1}(\delta)$. Similarly to the original safety bounds for randomized smoothing proposed by Cohen et al. (2019a), Thm. 4.3 does not require specific assumptions on the inner workings of f , nor the knowledge of its Lipschitz constant. Note that Thm. 4.3 is very general, since smoothing distributions such as $p^\# \mathcal{N}(0, \sigma^2 I)$, for different choices of p , allow one to sample the noise from a broad class of distributions. By choosing p wisely, we can sample from smoothing distributions that are appropriate for white-box multi-attacks as in Def. 3.1. Suitable choices of p may depend on the specific adversarial attacks considered. The proof of Thm. 4.4 provides an explicit choice of p , which is suitable for our case. Note that it is unclear what the relationship between the quantity $\|p^{-1}(\delta)\|_2$ and any known metric on the sample space x is, for a generic p . However, we show in the remainder of this section that it is possible to derive bounds for the Mahalanobis distance and the ℓ_2 norm using Thm. 4.3, for specific choices of p .

Probabilistic guarantees for the Mahalanobis distance. We first prove the following technical result, which allows us to build a suitable function p to apply Thm. 4.3 to our case.

Theorem 4.4. *Consider two random variables $X \sim \mathcal{N}(x; 0, \sigma^2 I)$ and $Y \sim \mathcal{N}(y; \mu, \Sigma)$. Suppose that Σ and I have the same dimensions. Furthermore, suppose that $\det(\sigma^2 I) = \det(\Sigma)$. Then, there exists a deterministic invertible function p such that:*

1. $Y = p^\# X$;
2. $\sqrt{(y - \mu)^T \Sigma^{-1} (y - \mu)} = \frac{1}{\sigma} \|p^{-1}(y)\|_2$ for all y in the support of Y .

Here, the function p is explicitly defined as $p(x) := \frac{1}{\sigma} Lx + \mu$, where L be a lower-triangular matrix that gives the Cholesky decomposition of Σ .

The proof of this theorem is deferred to App. C. By Thm. 4.4, we can apply Thm. 4.3 to the smoothing algorithm as in equation 3, to derive guarantees in terms of the Mahalanobis distance. The following lemma holds.

Lemma 4.5 (Probabilistic Guarantees for the Mahalanobis distance). *Consider a classifier f , and let $g(x)$ be the corresponding smoothed classifier as in equation 3. For a class $y_A \in Y$ suppose that*

there exist two constants $\underline{p}_A, \overline{p}_B \in [0, 1]$ such that

$$\mathbb{P}(f(x + \varepsilon) = y_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{y_B \neq y_A} \mathbb{P}(f(x + \varepsilon) = y_B),$$

with $\varepsilon \sim \mathcal{N}\left(0, \frac{\sigma^2}{\sqrt{\det(\Sigma)}} \Sigma\right)$. Then, it holds $g(\hat{x}) = y_A$ for all adversarial samples \hat{x} such that

$$\text{MAHL}(\hat{x} | x) \leq \frac{\sigma}{2 \sqrt[2d]{\det(\Sigma)}} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)).$$

This lemma allows one to derive probabilistic guarantees for randomized smoothing, in terms of the distance as in Def. 4.1. The proof of this result is deferred to App. C.

4.3 RELATIONSHIP WITH THE ℓ_2 NORM

In this section, we derive probabilistic guarantees for DSMOOTH, based on the ℓ_2 norm. There is a straightforward connection between the Mahalanobis distance and the ℓ_2 norm, as follows. For a matrix Σ as in Def. 4.1, denote with W any matrix such that $\Sigma = WW^T$. Then, it holds $\Sigma^{-1} = (W^{-1})^T W^{-1}$. Hence,

$$\text{MAHL}(\hat{x} | x) = \sqrt{(\hat{x} - x)^T (W^{-1})^T W^{-1} (\hat{x} - x)} = \|W^{-1}(\hat{x} - x)\|_2. \quad (4)$$

By equation 4, the Mahalanobis distance is the ℓ_2 norm after a *whitening transformation* (Kessy et al., 2018), i.e., a linear transformation that transforms a vector of random variables $\hat{x} - x$ with a known covariance matrix Σ into a set of new variables whose covariance is the identity matrix. In general, the matrix W in equation 4 is not uniquely defined. However, the resulting ℓ_2 norm $\|W^{-1}(\hat{x} - x)\|_2$ is equivalent across all these transformations, although some forms of W may have practical advantages over others (see, e.g., (Kessy et al., 2018)). We discuss common choices of W in App. D. By combining equation 4 with Lemma 4.5, we derive probabilistic guarantees for the ℓ_2 norm as follows.

Corollary 4.6 (Probabilistic guarantees for the ℓ_2 norm). *Consider a classifier f , and $g(x)$ be the corresponding smoothed classifier as in equation 3. For a class $y_A \in Y$ suppose that there exist two constants $\underline{p}_A, \overline{p}_B \in [0, 1]$ such that*

$$\mathbb{P}(f(x + \varepsilon) = y_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{y_B \neq y_A} \mathbb{P}(f(x + \varepsilon) = y_B),$$

with $\varepsilon \sim \mathcal{N}\left(0, \frac{\sigma^2}{\sqrt{\det(\Sigma)}} \Sigma\right)$. Then, it holds $g(\hat{x}) = y_A$ for all adversarial samples \hat{x} such that

$$\|W^{-1}(\hat{x} - x)\|_2 \leq \frac{\sigma}{2 \sqrt[2d]{\det(\Sigma)}} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)),$$

where W is any matrix such that $\Sigma = WW^T$.

We remark that, in general, the properties of W in Cor. 4.5 depend on the specific covariance matrix Σ . However, if the perturbations ε are sampled from an isotopic Gaussian distribution as in Cohen et al. (2019a), i.e., $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$, then Cor. 4.5 gives the same approximation guarantees as in Cohen et al. (2019a). In fact, consider a DSMOOTH algorithm with $\Sigma = \sigma^2 I$. For this algorithm, we can choose $W^{-1} = \frac{1}{\sigma} I$, and have that

$$\sqrt[2d]{\det(\Sigma)} = \sqrt[2d]{\det(\sigma^2 I)} = \sigma^2 \quad \text{and} \quad \|W^{-1}(\hat{x} - x)\|_2 = \frac{1}{\sigma} \|\hat{x} - x\|_2. \quad (5)$$

By substituting equation 5 in Cor. 4.5 we derive the same approximation guarantees as in Theorem 1 by Cohen et al. (2019a), which we restate for convenience.

Corollary 4.7 (Probabilistic guarantees for isotopic Gaussian noise, equivalent to Theorem 1 by Cohen et al. (2019a)). *Consider a classifier f , and $g(x)$ be the corresponding smoothed classifier as in equation 3, with $\Sigma = \sigma^2 I$. For a class $y_A \in Y$ suppose that there exist two constants $\underline{p}_A, \overline{p}_B \in [0, 1]$ such that*

$$\mathbb{P}(f(x + \varepsilon) = y_A) \geq \underline{p}_A \geq \overline{p}_B \geq \max_{y_B \neq y_A} \mathbb{P}(f(x + \varepsilon) = y_B),$$

with $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$. Then, it holds $g(\hat{x}) = y_A$ for all adversarial samples \hat{x} such that

$$\|\hat{x} - x\|_2 \leq \frac{\sigma}{2} (\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B)).$$

We remark that the bounds of Cor. 4.7 are known to be tight for isotopic Gaussian noise (Cohen et al., 2019a).

Table 1: Base classifiers and execution time of DSMOOTH on CIFAR-10. In this table, **Params.** (in millions) denotes the number of parameters, and **Time** (in seconds) denotes the average execution time required for a model to certify a datapoint, using the DCERT algorithm as in App. A. The execution time of DSMOOTH is similar to that of RANDSMOOTH and LSMOOTH on these base classifiers (see Table 3-5 in App. E.1).

Model	Params. (m)	Time (s)	Model	Params. (m)	Time (s)
resnet20	0.27	4.35 ± 0.08	mobilenetv2_x0_5	0.7	7.63 ± 0.33
resnet32	0.47	5.87 ± 0.24	mobilenetv2_x1_4	4.33	17.58 ± 0.54
resnet44	0.66	6.81 ± 0.09	shufflenetv2_x1_0	1.26	6.89 ± 0.09
resnet56	0.86	7.9 ± 0.05	shufflenetv2_x0_5	0.35	4.47 ± 0.1
vgg13_bn	9.94	7.43 ± 0.1	shufflenetv2_x2_0	5.37	13.73 ± 0.49
vgg16_bn	15.25	10.46 ± 0.21	repvgg_a0	7.84	18.74 ± 0.72
vgg19_bn	20.57	11.01 ± 0.09	repvgg_a1	12.82	26.28 ± 0.14
mobilenetv2_x1_0	2.24	12.2 ± 0.32	repvgg_a2	26.82	27.65 ± 0.26

Table 2: Base classifiers and execution time of DSMOOTH on IMAGENET. In this table, **Params.** (in millions) denotes the number of parameters, and **Time** (in seconds) denotes the average execution time required for a model to certify a datapoint, using the DCERT algorithm as in App. A. The execution time of DSMOOTH is similar to that of RANDSMOOTH and LSMOOTH on these base classifiers (see Table 4-6 in App. E.1).

Model	Params. (m)	Time (s)	Model	Params. (m)	Time (s)
resnet50	25.56	26.73 ± 0.21	wide_resnet50_2	68.88	40.24 ± 0.78
resnet152	60.19	107.52 ± 8.04	wide_resnet101_2	126.89	86.91 ± 11.97

5 EXPERIMENTS

The overall aim of the experiments is to demonstrate that DSMOOTH achieves good *certified accuracy* compared to baselines on complex adversarial attacks as in Def. 3.1. The certified accuracy is defined as the fraction of the test set, which a smoothed algorithm classifies correctly with a prediction that is certifiably robust within a ball of a given radius. Since DSMOOTH is a randomized smoothing classifier, it is not possible to compute this quantity exactly. Instead, we report on the approximate certified test set accuracy following previous related work, e.g., Cohen et al. (2019a). In addition to evaluating the certified accuracy, we also report on the execution time of DSMOOTH, and its sensitivity to different choices of the parameter k for the k -rank approximation as in Sec. 4.1.

In all the experiments we consider the SQUARE + FGSM multi-attack as detailed in Sec. 4.1, obtained as a combination of the Square Attack algorithm (Andriushchenko et al., 2020) and FGSM (Goodfellow et al., 2015). This attack applies a Square Attack to an input image (using ℓ_∞ norm), and then it applies a FGSM attack to the resulting adversarial sample (using ℓ_2 norm). In this experiment we use Square Attack with maximum perturbation 0.5 and 5000 queries. The FGSM attack uses maximum perturbation parameter 0.5. In App. F we show examples of CIFAR-10 (Fig. 6) and IMAGENET (Fig. 7) images corrupted with the smoothing distribution as in equation 3 for this type of attack.

5.1 OVERALL SET-UP

Base classifiers training. We consider various pre-trained classifiers, that achieve high accuracy on CIFAR-10 and IMAGENET respectively (see Table 1-2). We then fine-tune these classifiers to improve the robustness to adversarial attacks to SQUARE + FGSM as detailed in Sec. 3.1. Pre-trained models on CIFAR-10 (Table 1) are downloaded from <https://github.com/cheneafo/pytorch-cifar-models>, and pre-trained models on IMAGENET (Table 2) are downloaded from <https://github.com/pytorch/pytorch>. Fine-tuning consists of adjusting these models to a dataset that contains both CIFAR-10 training images and adversarial examples. The ratio of natural and adversarial examples is 50 : 50. In this work, we opt for a simple training procedure, to highlight the benefits of our method against baselines. However, we believe that the certified accuracy of our method could be further improved by considering more complex adversarial training procedures, such as Wong & Kolter (2021).

Baselines. We compared our smoothing algorithm in equation 3 to two baseline approaches for certified robustness: the standard randomized smoothing algorithm by (Cohen et al., 2019a) (RANDSMOOTH), and the approach by Teng et al. (2020) (LSMOOTH). The randomized smoothing algorithm by Cohen et al. (2019a) provides certification guarantees in terms of the ℓ_2 norm, whereas the algorithm by Teng et al. (2020) provides guarantees in terms of the ℓ_1 norm. We do not consider randomized smoothing techniques with certification guarantees in terms of ℓ_p norms with $p > 2$, since impossibility results are known for increasing p (Yang et al., 2020; Blum et al., 2020; Kumar et al., 2020). Specifically, we do not consider any certification mechanism for the ℓ_∞ norm, since the isotropic Gaussian distribution as in RANDSMOOTH is optimal for defending against ℓ_∞ attacks, if we don't use a more powerful technique than Neyman-Pearson (Yang et al., 2020).

System. The system used features multiple Intel® Xeon® Gold 6252 CPUs, each with a base clock speed of 2.10 GHz, operating at various frequencies between 2011 MHz and 2800 MHz. The system also includes six NVIDIA GPUs for more intensive graphics and computational workloads. These are two NVIDIA GPUs with a 64-bit width and clock speed of 33 MHz, and four NVIDIA GPUs of the GV102 model with a 64-bit width, operating at a clock speed of 33 MHz.

5.2 RESULTS ON CIFAR-10

Execution time. We test the performance of DSMOOTH. To this end, we run the DCERT algorithm, as detailed in App. A, on various base models. Parameters for DCERT are $\alpha = 0.001$, $n_0 = 100$ Monte Carlo samples for selection and $n = 100000$ samples for estimation. With this parameters choice, there is at most 0.001 probability that DCERT returns a radius that is not robust (see App. A). For each base classifier, we test our algorithm on 500 images from CIFAR-10 and we report on the average execution time (in seconds) to certify a single image. The results are reported in Table 1. Overall we observe that DSMOOTH is scalable to all models, although for models with several million parameters, such as RepVGG_a2, the performance decreases. We remark that the performance of our algorithm is similar to the performance of previous algorithms, e.g., the algorithms by Cohen et al. (2019a); Teng et al. (2020). We refer the reader to Table 3-5 in App. E.1 for the execution time of previous algorithms.

Comparison against baselines. We run the DCERT algorithm (App. A) against baselines with parameters $\alpha = 0.001$, $n_0 = 100$ samples for selection and $n = 100000$ samples for estimation. For each base classifier in Table 1, we test DCERT and baselines on 500 images from CIFAR-10. The results are displayed in Fig. 1, where we observe that in all cases our algorithm performs significantly better than the baselines. These results demonstrate that DSMOOTH is suitable to handle complex adversarial attacks as in Def. 3.1, whereas RANDSMOOTH and LSMOOTH are unsuitable to that end. In fact, in most cases the certified accuracy of RANDSMOOTH and LSMOOTH is approximately 0.1. Since CIFAR-10 has only 10 classes, these results suggest that RANDSMOOTH and LSMOOTH do not perform significantly better than uniform random sampling.

Additional experiments. In App. E.2 we provide additional experiments on CIFAR-10 to determine the effect of different choices of α and number of samples for selection n on the performance of DCERT.

5.3 RESULTS ON IMAGENET

Execution time. We evaluate the effectiveness of our smoothing algorithm as described in equation 3. To achieve this, we apply the DCERT algorithm, as outlined in App. A, across different base models. For DCERT, we use parameters $\alpha = 0.001$, $n_0 = 100$ Monte Carlo samples for selection, and $n = 1000$ samples for estimation. In this scenario, we approximate the matrix Σ from equation 2 using a PCA algorithm, as explained in Sec. 4.1, with a rank- k approximation where $k = 1000$. Each base classifier is tested on 500 images from IMAGENET, and we measure the average execution time (in seconds) required to certify a single image. The results are summarized in Table 2. Overall, DSMOOTH demonstrates scalability to very large models, and its performance is comparable to that of previous algorithms (see Table 3-5 in App. E.1).

Comparison against baselines. Once again, we evaluated our smoothing algorithm from equation 3 against baselines. We apply the DCERT algorithm (see App. A) with parameters $\alpha = 0.001$, $n_0 = 100$ samples for selection, and $n = 1000$ samples for estimation. Due to the large size of Σ , we use a rank- k approximation Σ_k with $k = 1000$. For each base classifier listed in Table 2, we evaluate DCERT and the baseline methods on 500 images from CIFAR-10. The results are presented in Fig. 2,

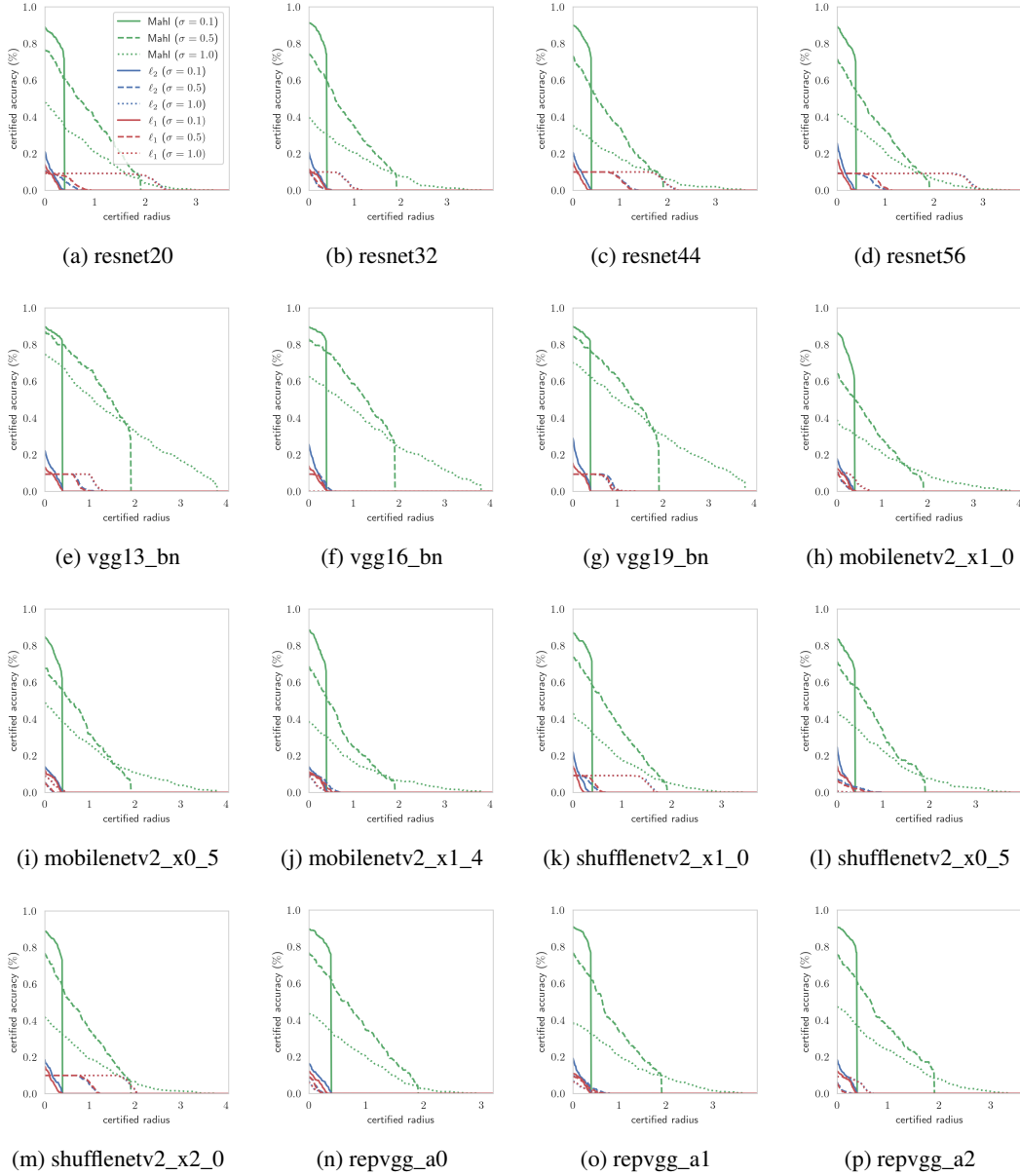


Figure 1: Approximate certified accuracy of DSMOOTH (MAHL in the legend), RANDSMOOTH (ℓ_2 in the legend) and LSMOOTH (ℓ_1 in the legend) on CIFAR-10 for various base models as in Table 1. DSMOOTH is significantly better than baselines.

showing that our algorithm consistently outperforms the baselines. As with the CIFAR-10 results, this demonstrates that DSMOOTH is effective against complex adversarial attacks as defined in Def. 3.1, while RANDSMOOTH and LSMOOTH are inadequate for this purpose.

Ablation study on the rank- k approximation of Σ . We conclude with a set of experiments to determine if our results are sensitive to the rank k of the PCA approximation of the covariance matrix Σ . To this end, we run the DCERT algorithm with the smoothing distribution as in equation 3, for $k = 10, 100, 1000, 10000$. Each run uses the parameters $\sigma = 0.5$, $\alpha = 0.001$, $n_0 = 100$ samples for selection and $n = 1000$ samples for estimation. The results are displayed in Fig. 3. The results suggest that DSMOOTH is not very sensitive to different choices of k .

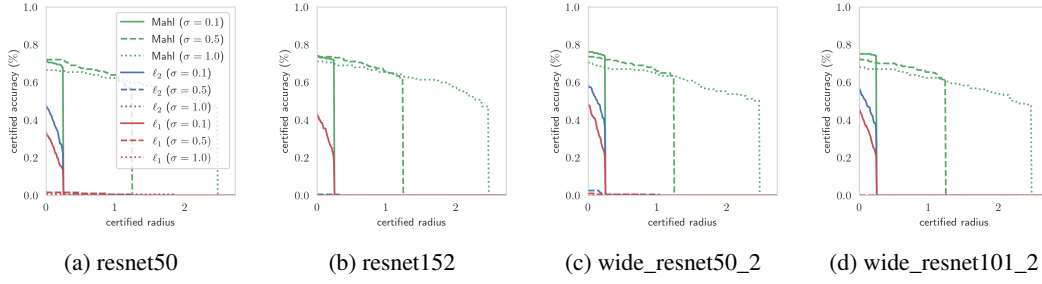


Figure 2: Approximate certified accuracy of DSMOOTH (MAHL in the legend), RANDSMOOTH (ℓ_2 in the legend) and LSMOOTH (ℓ_1 in the legend) on IMAGENET for various base models as in Table 2. We observe that DSMOOTH comes out on top.

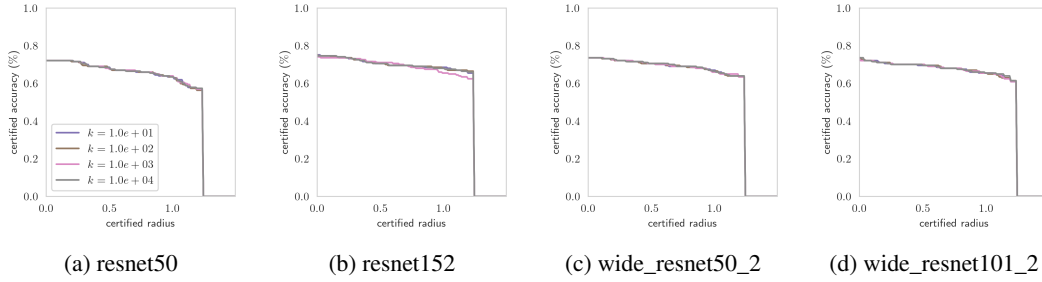


Figure 3: Approximate certified accuracy of DSMOOTH for different choices of k on IMAGENET, for various base models as in Table 2. We observe that the parameter k does not significantly affect the performance of DSMOOTH.

6 DISCUSSION

In this paper, we introduced a novel certification method based on randomized smoothing (see equation 3) to enhance the robustness of machine learning models against complex adversarial attacks, including combinations of multiple attack types (see Sec. 4.1). Our approach generalizes the existing framework of randomized smoothing by incorporating more flexible noise distributions, allowing for robustness guarantees across a wider range of adversarial threats, such as SQUARE+FGSM (see Sec. 4.1). Through extensive experiments on CIFAR-10 (see Sec. 5.2) and IMAGENET (see Sec. 5.3), we demonstrated that our method consistently outperforms state-of-the-art defenses in terms of certified accuracy (see Fig. 1-2).

However, much like previous work (see, e.g., Cohen et al. (2019a); Teng et al. (2020)), our proposed method still faces several limitations. The effectiveness of DSMOOTH is constrained by its reliance on Monte Carlo sampling, which can be computationally expensive on very large models. Additionally, while our approach extends robustness beyond the standard ℓ_2 norm, it may not yet fully capture the complexities of all possible adversarial threats.

Future work could address these limitations by developing more efficient sampling techniques, or by leveraging neural architecture search to identify base classifiers that are inherently more robust to adversarial perturbations. Furthermore, exploring alternative noise distributions and adaptive smoothing strategies could further enhance robustness against a broader array of adversarial threats.

REFERENCES

- Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4):433–459, 2010.
- Alaa Anani, Tobias Lorenz, Bernt Schiele, and Mario Fritz. Adaptive hierarchical certification for segmentation using randomized smoothing. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.

- Brendon G. Anderson and Somayeh Sojoudi. Certified robustness via locally biased randomized smoothing. In *Learning for Dynamics and Control Conference*, pp. 207–220. PMLR, 2022.
- Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: A query-efficient black-box adversarial attack via random search. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXIII*, volume 12368 of *Lecture Notes in Computer Science*, pp. 484–501. Springer, 2020.
- Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 274–283. PMLR, 2018.
- Avrim Blum, Travis Dick, Naren Manoj, and Hongyang Zhang. Random smoothing might be unable to certify ∞ robustness for high-dimensional images. *J. Mach. Learn. Res.*, 21:211:1–211:21, 2020.
- Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. Efficient robustness certificates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1003–1013. PMLR, 2020. URL <http://proceedings.mlr.press/v119/bojchevski20a.html>.
- Tom B. Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *CoRR*, abs/1712.09665, 2017a.
- Tom B. Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *arXiv preprint arXiv:1712.09665*, 2017b.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. *2017 IEEE symposium on security and privacy (sp)*, 2017a.
- Nicholas Carlini and David A. Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In Bhavani Thuraisingham, Battista Biggio, David Mandell Freeman, Brad Miller, and Arunesh Sinha (eds.), *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISec@CCS 2017, Dallas, TX, USA, November 3, 2017*, pp. 3–14. ACM, 2017b.
- Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pp. 39–57. IEEE Computer Society, 2017c.
- Jeremy Cohen, Elan Rosenfeld, and J. Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1310–1320. PMLR, 2019a.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pp. 1310–1320. PMLR, 2019b.
- Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International Conference on Machine Learning*. PMLR, 2020.
- Krishnamurthy Dvijotham, Robert Stanforth, Sven Gowal, Timothy A. Mann, and Pushmeet Kohli. A dual approach to scalable verification of deep networks. In Amir Globerson and Ricardo Silva (eds.), *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence, UAI 2018, Monterey, California, USA, August 6-10, 2018*, pp. 550–559. AUAI Press, 2018.
- Rüdiger Ehlers. Formal verification of piece-wise linear feed-forward neural networks. In Deepak D’Souza and K. Narayan Kumar (eds.), *Automated Technology for Verification and Analysis - 15th International Symposium, ATVA 2017, Pune, India, October 3-6, 2017, Proceedings*, volume 10482 of *Lecture Notes in Computer Science*, pp. 269–286. Springer, 2017.

- Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring the landscape of spatial robustness. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1802–1811. PMLR, 2019.
- Marc Fischer, Maximilian Baader, and Martin T. Vechev. Certified defense to image transformations via randomized smoothing. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6573*, 2014.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- Shixiang Gu and Luca Rigazio. Towards deep neural network architectures robust to adversarial examples. *arXiv preprint arXiv:1412.5068*, 2014.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 770–778. IEEE Computer Society, 2016.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *CoRR*, abs/1704.04861, 2017.
- J. Edward Hu, Adith Swaminathan, Hadi Salman, and Greg Yang. Improved image wasserstein attacks and defenses. *CoRR*, abs/2004.12478, 2020.
- Xiaowei Huang, Marta Kwiatkowska, Sen Wang, and Min Wu. Safety verification of deep neural networks. In Rupak Majumdar and Viktor Kuncak (eds.), *Computer Aided Verification - 29th International Conference, CAV 2017, Heidelberg, Germany, July 24-28, 2017, Proceedings, Part I*, volume 10426 of *Lecture Notes in Computer Science*, pp. 3–29. Springer, 2017.
- Kenneth Hung and William Fithian. Rank verification for exponential families. *the annals of statistics. The Annals of Statistics*, 2(04):758–782, 2019.
- Gaojie Jin, Xinpeng Yi, Dengyu Wu, Ronghui Mu, and Xiaowei Huang. Randomized adversarial training via taylor expansion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pp. 16447–16457. IEEE, 2023.
- Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Reluplex: An efficient SMT solver for verifying deep neural networks. In Rupak Majumdar and Viktor Kuncak (eds.), *Computer Aided Verification - 29th International Conference, CAV 2017, Heidelberg, Germany, July 24-28, 2017, Proceedings, Part I*, volume 10426 of *Lecture Notes in Computer Science*, pp. 97–117. Springer, 2017.
- Agnan Kessy, Alex Lewin, and Korbinian Strimmer. Optimal whitening and decorrelation. *The American Statistician*, 72(4):309–314, 2018.
- Klim Kireev, Maksym Andriushchenko, and Nicolas Flammarion. On the effectiveness of adversarial training against common corruptions. In James Cussens and Kun Zhang (eds.), *Uncertainty in Artificial Intelligence, Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence, UAI 2022, 1-5 August 2022, Eindhoven, The Netherlands*, volume 180 of *Proceedings of Machine Learning Research*, pp. 1012–1021. PMLR, 2022.

- Aounon Kumar, Alexander Levine, Tom Goldstein, and Soheil Feizi. Curse of dimensionality on randomized smoothing for certifiable robustness. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5458–5467. PMLR, 2020.
- Anupriya Kumari, Devansh Bhardwaj, Sukrit Jindal, and Sarthak Gupta. Trust, but verify: A survey of randomized smoothing techniques. *CoRR*, abs/2312.12608, 2023.
- Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. 2016.
- Marta Kwiatkowska and Xiyue Zhang. When to trust AI: advances and challenges for certification of neural networks. In Maria Ganzha, Leszek A. Maciaszek, Marcin Paprzycki, and Dominik Slezak (eds.), *Proceedings of the 18th Conference on Computer Science and Intelligence Systems, FedCSIS 2023, Warsaw, Poland, September 17-20, 2023*, volume 35 of *Annals of Computer Science and Information Systems*, pp. 25–37, 2023.
- Cassidy Laidlaw and Soheil Feizi. Functional adversarial attacks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 10408–10418, 2019.
- Cassidy Laidlaw, Sahil Singla, and Soheil Feizi. Perceptual adversarial robustness: Defense against unseen threat models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE symposium on security and privacy (SP)*, pp. 656–672. IEEE, 2019.
- Guang-He Lee, Yang Yuan, Shiyu Chang, and Tommi Jaakkola. Tight certificates of adversarial robustness for randomly smoothed classifiers. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Alexander Levine and Soheil Feizi. Wasserstein smoothing: Certified robustness against wasserstein adversarial attacks. In Silvia Chiappa and Roberto Calandra (eds.), *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online [Palermo, Sicily, Italy]*, volume 108 of *Proceedings of Machine Learning Research*, pp. 3938–3947. PMLR, 2020a.
- Alexander Levine and Soheil Feizi. (de) randomized smoothing for certifiable defense against patch attacks. *Advances in Neural Information Processing Systems*, 33:6465–6475, 2020b.
- Alexander Levine and Soheil Feizi. Robustness certificates for sparse adversarial attacks by randomized ablation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 2020c.
- Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive noise. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yujie Li, Xing Xu, Jinhui Xiao, Siyuan Li, and Heng Tao Shen. Adaptive square attack: Fooling autonomous cars with adversarial traffic signs. *IEEE Internet Things J.*, 8(8):6337–6347, 2021.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017a.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017b.

- Yuhao Mao, Mark Niklas Müller, Marc Fischer, and Martin T. Vechev. Connecting certified and adversarial training. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Yuhao Mao, Mark Niklas Müller, Marc Fischer, and Martin T. Vechev. Understanding certified training with interval bound propagation. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- David J. Miller, Zhen Xiang, and George Kesidis. Adversarial learning targeting deep neural network classification: A comprehensive review of defenses against attacks. *Proceedings of the IEEE*, 108(3):402–433, 2020.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 2574–2582. IEEE Computer Society, 2016.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *2016 IEEE symposium on security and privacy (SP)*, pp. 582–597. IEEE, 2016a.
- Nicolas Papernot, Patrick D. McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *IEEE European Symposium on Security and Privacy, EuroS&P 2016, Saarbrücken, Germany, March 21-24, 2016*, pp. 372–387. IEEE, 2016b.
- Edward Raff, Jon Barker, Jared Sylvester, Robert Brandon, Bryan Catanzaro, and Charles K. Nicholas. Malware detection by eating a whole exe. In *Workshops at the thirty-second AAAI conference on artificial intelligence*, 2018.
- Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- Yan Scholten, Jan Schuchardt, Aleksandar Bojchevski, and Stephan Günnemann. Hierarchical randomized smoothing. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Erfan Shayegani, Md Abdullah Al Mamun, Yu Fu, Pedram Zaree, Yue Dong, and Nael B. Abu-Ghazaleh. Survey of vulnerabilities in large language models revealed by adversarial attacks. *CoRR*, abs/2310.10844, 2023.
- 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.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- Jiaye Teng, Guang-He Lee, and Yang Yuan. ℓ_1 adversarial robustness certificates: a randomized smoothing approach. 2020.

- Florian Tramèr, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. In Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Hong Wang, Yuefan Deng, Shinjae Yoo, Haibin Ling, and Yuewei Lin. AGKD-BML: defense against adversarial attack by attention guided knowledge distillation and bi-directional metric learning. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 7638–7647. IEEE, 2021.
- Fabian Woitschek and Georg Schneider. Physical adversarial attacks on deep neural networks for traffic sign recognition: A feasibility study. *CoRR*, abs/2302.13570, 2023.
- Eric Wong and J. Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 5283–5292. PMLR, 2018.
- Eric Wong and J. Zico Kolter. Learning perturbation sets for robust machine learning. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- Eric Wong, Frank R. Schmidt, and J. Zico Kolter. Wasserstein adversarial examples via projected sinkhorn iterations. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6808–6817. PMLR, 2019.
- Greg Yang, Tony Duan, J. Edward Hu, Hadi Salman, Ilya Razenshteyn, and Jerry Li. Randomized smoothing of all shapes and sizes. In *International Conference on Machine Learning*, pp. 10693–10705. PMLR, 2020.
- Jiawei Zhang, Linyi Li, Ce Zhang, and Bo Li. Care: Certifiably robust learning with reasoning via variational inference. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pp. 554–574. IEEE, 2023.
- Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Cheung, and Min Lin. On evaluating adversarial robustness of large vision-language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *CoRR*, abs/2307.15043, 2023.

Algorithm 1 Pseudocode for Certification and Prediction

```

//evaluate  $g$  at  $x$ 
function DPRED ( $f, \Sigma, \sigma, x, n, \alpha$ ) :
    for  $i \in [n]$  do  $\hat{c}_i \leftarrow f(x + \varepsilon)$ , with  $\varepsilon \sim \mathcal{N}\left(0, \frac{\sigma^2}{\sqrt{d \det(\Sigma)}} \Sigma\right)$ ;
     $\hat{c}_A, \hat{c}_B \leftarrow$  top two indices in  $\{\hat{c}_1, \dots, \hat{c}_n\}$ ;
     $n_A, n_B \leftarrow$  frequency of  $\hat{c}_A, \hat{c}_B$  in  $\{\hat{c}_1, \dots, \hat{c}_n\}$ ;
    if BINOMPVALUE( $n_A, n_A + n_B, 0.5$ )  $\leq \alpha$  then return  $\hat{c}_A$ ;
    else return ABSTAIN;

//certify the robustness of  $g$  around  $x$ 
function DCERT ( $f, p_\theta, \sigma, x, n_0, n, \alpha$ ) :
    for  $i \in [n_0]$  do  $\hat{c}_i^0 \leftarrow f(x + \varepsilon)$ , with  $\varepsilon \sim \mathcal{N}\left(0, \frac{\sigma^2}{\sqrt{d \det(\Sigma)}} \Sigma\right)$ ;

    for  $i \in [n]$  do  $\hat{c}_i \leftarrow f(x + \varepsilon)$ , with  $\varepsilon \sim \mathcal{N}\left(0, \frac{\sigma^2}{\sqrt{d \det(\Sigma)}} \Sigma\right)$ ;

     $\hat{c}_A^0 \leftarrow$  top index in  $\{\hat{c}_1^0, \dots, \hat{c}_{n_0}^0\}$ ;
     $n_A \leftarrow$  frequency of  $\hat{c}_A^0$  in  $\{\hat{c}_1, \dots, \hat{c}_n\}$ ;
     $p_A \leftarrow$  LOWERCONFBOUND( $n_A, n, 1 - \alpha$ );
    if  $p_A > 0.5$  return prediction  $\hat{c}_A^0$  and radius  $\sigma \Phi^{-1}(p_A)$ ;
    else return ABSTAIN;

```

A PRACTICAL ALGORITHMS

Following previous work (Cohen et al., 2019a), we present practical Monte Carlo algorithms for evaluating the smoothed classifier $g(x)$ and certifying the robustness of g around x . These algorithms are called DPRED and DCERT respectively. Pseudocode for these procedures is presented in Alg. 1.

DPRED. To evaluate the smoothed classifier’s prediction $g(x)$, one must identify the class c_A that has the highest weight in the categorical distribution $f(p(x, e))$ with $e \sim \mathcal{N}(0, \sigma^2 I)$. To estimate c_A we propose DPRED, which follows the standard approach outlined by (Cohen et al., 2019a). Pseudocode for this algorithm is presented in Alg. 1, where the function BINOMPVALUE($n_A, n_A + n_B, 0.5$) returns the p -value of the two-sided hypothesis test that $n_A \sim \text{BINOMIAL}(n_A + n_B, p)$. Intuitively, the function DPRED involves drawing n samples from $f \circ p(x, e)$. The class c_A that appears most frequently is noted. If c_A significantly outnumbers other classes, DPRED returns c_A ; otherwise, it abstains from making a prediction. Following Cohen et al. (2019a), the abstention threshold is calibrated using the hypothesis test from Hung & Fithian (2019) to ensure that the probability of an incorrect prediction is limited to α , thus guaranteeing a controlled error rate. In fact, from Proposition 1 by Cohen et al. (2019a) it follows that the probability that DPRED returns a class other than $g(x)$ is at most α . The following lemma holds.

Lemma A.1 (Following Proposition 1 by Cohen et al. (2019a)). *With probability at least $1 - \alpha$, DPRED either abstain or it returns $g(x)$.*

For a proof of this result we refer the reader to Appendix C in Cohen et al. (2019a).

DCERT. We provide an algorithm for certifying the robustness of g around an input x , which we refer to as DCERT. Again, this algorithm follows the standard certification procedure for randomized smoothing by Cohen et al. (2019a). Pseudocode for DCERT is presented in Alg. 1, where the function LOWERCONFBOUND($n_A, n, 1 - \alpha$) returns a one-sided $1 - \alpha$ lower confidence interval for the Binomial parameter p given a sample $k \sim \text{BINOMIAL}(n, p)$. Intuitively, DCERT uses a small number of samples n_0 from $f(p(x, e))$ to identify c_A . Then, it uses a large number of samples n to estimate p_A . Finally, DCERT sets $p_B = 1 - p_A$. The following lemma. Following Proposition 2 by Cohen et al. (2019a), we have that with probability at least $1 - \alpha$ over the randomness in DCERT, if DCERT returns a class c_A and a radius R (i.e. does not abstain), then we have the robustness guarantee $g(\hat{x}) = c_A$ whenever $\text{MAHL}(\hat{x} | x) \leq R$. Formally, the following lemma holds.

Lemma A.2 (Following Proposition 2 by Cohen et al. (2019a)). *With probability at least $1 - \alpha$, if DCERT does not abstain, then g predicts c_A within radius R around x , i.e., $g(\hat{x}) = c_A$ for all adversarial examples \hat{x} such that $\text{MAHL}(\hat{x} \mid x) \leq R$.*

For a proof of this result we refer the reader to Appendix C in Cohen et al. (2019a).

B PROOF OF THEOREM 4.3

B.1 PRELIMINARY RESULTS

In order to prove this theorem, we first consider the following auxiliary results.

Lemma B.1 (Neyman-Pearson, following Lemma 3 by Cohen et al. (2019a)). *Let X and Y be random variables in \mathbb{R}^d with the same support. Let $\ell: \mathbb{R}^d \rightarrow \{0, 1\}$ be any random or deterministic function. Then, it holds:*

- *If $S = \{z \in \mathbb{R}^d: \mathbb{P}(X = z) \leq t\mathbb{P}(Y = z)\}$ for some $t \geq 0$ and $\mathbb{P}(\ell(X) = 1) \geq \mathbb{P}(X \in S)$, then $\mathbb{P}(\ell(Y) = 1) \geq \mathbb{P}(X \in S)$;*
- *If $S = \{z \in \mathbb{R}^d: \mathbb{P}(Y = z) \geq t\mathbb{P}(X = z)\}$ for some $t \geq 0$ and $\mathbb{P}(\ell(X) = 1) \leq \mathbb{P}(X \in S)$, then $\mathbb{P}(\ell(Y) = 1) \leq \mathbb{P}(X \in S)$.*

We also consider the following auxiliary lemma.

Lemma B.2 (Neyman-Pearson for the Push-Forward Measure). *Consider an invertible measurable mapping $p: \Omega \subseteq \mathbb{R}^k \rightarrow \mathbb{R}^d$ and define the random variables*

$$\begin{aligned}\hat{X}_0 &:= p(e) \text{ and } E_0 := e, \text{ with } e \sim \mathcal{N}(e; 0, \sigma^2 I), \\ \hat{X}_\delta &:= p(e) \text{ and } E_\delta := e, \text{ with } e \sim \mathcal{N}(e; \delta, \sigma^2 I).\end{aligned}$$

Let $\ell: \mathbb{R}^d \rightarrow \{0, 1\}$ be any deterministic or random function. The following statements hold.

- *If $S := \{z \in \mathbb{R}^d: \delta^T p^{-1}(z) \leq \beta\}$ for a constant β , and $\mathbb{P}(\ell(\hat{X}_0) = 1) \geq \mathbb{P}(\hat{X}_0 \in S)$, then $\mathbb{P}(\ell(\hat{X}_\delta) = 1) \geq \mathbb{P}(\hat{X}_\delta \in S)$.*
- *If $S := \{z \in \mathbb{R}^d: \delta^T p^{-1}(z) \geq \beta\}$ for a constant β , and $\mathbb{P}(\ell(\hat{X}_0) = 1) \leq \mathbb{P}(\hat{X}_0 \in S)$, then $\mathbb{P}(\ell(\hat{X}_\delta) = 1) \leq \mathbb{P}(\hat{X}_\delta \in S)$.*

Proof. To simplify the notation, we denote with $p_i^{-1}(z)$ the i -th component of the inverse map $p^{-1}(z)$. We further define

$$\mu_0(e) := \mathcal{N}(e; 0, \sigma^2 I) \quad \text{and} \quad \mu_\delta(e) := \mathcal{N}(e; \delta, \sigma^2 I).$$

We further denote with $p^\# \mu_0(z) := \mu_0(p^{-1}(z))$ and $p^\# \mu_\delta(z) := \mu_\delta(p^{-1}(z))$ the corresponding push-forward measures. By Lemma B.1 it suffices to show that for any β , there exists a constant $t > 0$ for which

$$\{z \in \mathbb{R}^d: \delta^T p^{-1}(z) \leq \beta\} = \left\{z \in \mathbb{R}^d: \frac{p^\# \mu_\delta(z)}{p^\# \mu_0(z)} \leq t\right\} \quad (6)$$

and

$$\{z \in \mathbb{R}^d: \delta^T p^{-1}(z) \geq \beta\} = \left\{z \in \mathbb{R}^d: \frac{p^\# \mu_\delta(z)}{p^\# \mu_0(z)} \geq t\right\}. \quad (7)$$

To this end, we compute the likelihood ratio for the ratio of the chosen push-forward measures as

$$\frac{p^\# \mu_\delta(z)}{p^\# \mu_0(z)} = \frac{\exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^d (p_i^{-1}(z) - \delta_i)^2\right\}}{\exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^d (p_i^{-1}(z))^2\right\}} \quad (8)$$

$$= \exp\left\{\frac{1}{2\sigma^2} \sum_{i=1}^d (\delta_i p_i^{-1}(z) - \delta_i^2)\right\} \quad (9)$$

$$= \exp\{a\delta^T p^{-1}(z) - b\}, \quad (10)$$

where Eq. equation 8 follows from the definition of μ_0 and μ_δ , Eq. equation 9 follows from standard calculations, and Eq. equation 10 follows by choosing $a := 1/\sigma^2$ and $b := \sum_{i=1}^d \delta_i^2/\sigma^2$. Therefore, for any constant β , we may take $t = \exp\{a\beta + b\}$, noticing that

$$\begin{aligned}\delta^T p^{-1}(z) \leq \beta &\Leftrightarrow \exp\{\alpha \delta^T p^{-1}(z) - \beta\} \leq t, \\ \delta^T p^{-1}(z) \geq \beta &\Leftrightarrow \exp\{\alpha \delta^T p^{-1}(z) - \beta\} \geq t.\end{aligned}$$

Then, Eq. equation 6 and Eq. equation 7 hold, and the claim follows. \square

B.2 PROOF OF THE MAIN RESULT

The proof of Theorem 4.3 relies on a few additional simple lemmas, which we present before giving the main proof.

Lemma B.3. *Define the random variables*

$$\begin{aligned}\hat{X}_0 &:= p_\theta(x, e) \text{ and } E_0 := e, \text{ with } e \sim \mathcal{N}(e; 0, \sigma^2 I), \\ Z &:= e, \text{ with } e \sim \mathcal{N}(e; 0, I),\end{aligned}$$

and consider the half spaces

$$\begin{aligned}S_{\leq} &:= \{\delta^T p_x^{-1}(\hat{x}) \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)\}, \\ S_{\geq} &:= \{\delta^T p_x^{-1}(\hat{x}) \geq \sigma \|\delta\| \Phi^{-1}(1 - \overline{p}_B)\}.\end{aligned}$$

Then, it holds $\mathbb{P}(\hat{X}_0 \in S_{\leq}) = \underline{p}_A$ and $\mathbb{P}(\hat{X}_0 \in S_{\geq}) = \overline{p}_B$.

Proof. To show the first part of the claim, note that it holds

$$\begin{aligned}\mathbb{P}(\hat{X}_0 \in S_{\leq}) &= \mathbb{P}(\delta^T p_x^{-1}(\hat{X}_0) \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)) \\ &= \mathbb{P}(\delta^T E_0 \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)) \\ &= \mathbb{P}(\sigma \|\delta\| Z \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)) \\ &= \Phi(\Phi^{-1}(\underline{p}_A)) \\ &= \underline{p}_A\end{aligned}$$

We prove the second part of the claim in a similar fashion. It holds

$$\begin{aligned}\mathbb{P}(\hat{X}_0 \in S_{\geq}) &= \mathbb{P}(\delta^T p_x^{-1}(\hat{X}_0) \geq \sigma \|\delta\| \Phi^{-1}(\overline{p}_B)) \\ &= \mathbb{P}(\delta^T E_0 \geq \sigma \|\delta\| \Phi^{-1}(\overline{p}_B)) \\ &= \mathbb{P}(\sigma \|\delta\| Z \geq \sigma \|\delta\| \Phi^{-1}(\overline{p}_B)) \\ &= \Phi(\Phi^{-1}(\overline{p}_B)) \\ &= \overline{p}_B\end{aligned}$$

\square

Lemma B.4. *Define the random variables*

$$\begin{aligned}\hat{X}_\delta &:= p_\theta(x, e) \text{ and } E_0 := e, \text{ with } e \sim \mathcal{N}(e; 0, \sigma^2 I), \\ Z &:= e, \text{ with } e \sim \mathcal{N}(e; 0, I),\end{aligned}$$

and consider the half spaces

$$\begin{aligned}S_{\leq} &:= \{\delta^T p_x^{-1}(\hat{x}) \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)\}, \\ S_{\geq} &:= \{\delta^T p_x^{-1}(\hat{x}) \geq \sigma \|\delta\| \Phi^{-1}(1 - \overline{p}_B)\}.\end{aligned}$$

Then, it holds $\mathbb{P}(\hat{X}_\delta \in S_{\leq}) = \Phi(\Phi^{-1}(\underline{p}_A) - \|\delta\|/\sigma)$ and $\mathbb{P}(\hat{X}_\delta \in S_{\geq}) = \Phi(\Phi^{-1}(\overline{p}_B) - \|\delta\|/\sigma)$.

Proof. To show the first part of the claim, note that it holds

$$\begin{aligned}
\mathbb{P}(\hat{X}_\delta \in S_{\leq}) &= \mathbb{P}(\delta^T p_x^{-1}(\hat{X}_\delta) \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)) \\
&= \mathbb{P}(\delta^T E_\delta \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)) \\
&= \mathbb{P}(\sigma \delta^T Z + \|\delta\|^2 \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)) \\
&= \mathbb{P}(Z \leq \Phi^{-1}(\underline{p}_A) - \|\delta\|/\sigma) \\
&= \Phi(Z \leq \Phi^{-1}(\underline{p}_A) - \|\delta\|/\sigma)
\end{aligned}$$

We prove the second part of the claim in a similar fashion. It holds

$$\begin{aligned}
\mathbb{P}(\hat{X}_\delta \in S_{\geq}) &= \mathbb{P}(\delta^T p_x^{-1}(\hat{X}_\delta) \geq \sigma \|\delta\| \Phi^{-1}(1 - \overline{p}_B)) \\
&= \mathbb{P}(\delta^T E_\delta \geq \sigma \|\delta\| \Phi^{-1}(1 - \overline{p}_B)) \\
&= \mathbb{P}(\sigma \delta^T Z + \|\delta\|^2 \leq \sigma \|\delta\| \Phi^{-1}(1 - \overline{p}_B)) \\
&= \mathbb{P}(Z \leq \Phi^{-1}(\overline{p}_B) - \|\delta\|/\sigma) \\
&= \Phi(Z \leq \Phi^{-1}(\overline{p}_B) - \|\delta\|/\sigma)
\end{aligned}$$

□

Lemma B.5. Define the random variables

$$\begin{aligned}
\hat{X}_0 &:= p_\theta(x, e) \text{ and } E_0 := e, \text{ with } e \sim \mathcal{N}(e; 0, \sigma^2 I), \\
\hat{X}_\delta &:= p_\theta(x, e) \text{ and } E_\delta := e, \text{ with } e \sim \mathcal{N}(e; \delta, \sigma^2 I), \\
Z &:= e, \text{ with } e \sim \mathcal{N}(e; 0, I),
\end{aligned}$$

and consider the half spaces

$$\begin{aligned}
S_{\leq} &:= \{\delta^T p_x^{-1}(\hat{x}) \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)\}, \\
S_{\geq} &:= \{\delta^T p_x^{-1}(\hat{x}) \geq \sigma \|\delta\| \Phi^{-1}(1 - \overline{p}_B)\}.
\end{aligned}$$

Then, it holds $\mathbb{P}(\hat{X}_\delta \in S_{\leq}) \geq \mathbb{P}(\hat{X}_\delta \in S_{\geq})$ if and only if $\|\delta\| \leq R$.

Proof. The proof of this claim is an application of Lemma B.4. In fact, algebra shows that the claim holds if and only if $\|\delta\| \leq \frac{\sigma}{2}(\Phi^{-1}(\underline{p}_A) - \Phi^{-1}(\overline{p}_B))$. □

We now have all the necessary tools to prove Thm. 4.3.

Proof of Thm. 4.3. Fix a constant $R > 0$, and suppose that it holds

$$(g \circ p_\theta)(x, e) = y_A \text{ for all } x \text{ and } \|e\|_2 < R. \quad (11)$$

Then, the claim holds. In fact, if the reconstruction loss of the decoder $p_\theta(x, e)$ is zero, it holds $\hat{x} = \mu_\phi(\hat{x}, x)$ from which it follows that $g(\hat{x}) = (g \circ p_\theta)(\mu_\phi(\hat{x}, x))$. Hence, from Eq. equation 11 it follows that $g(\hat{x}) = y_A$ for all \hat{x} such that $\|\mu_\phi(\hat{x}, x)\| \leq R$. The claim follows using Jensen's inequality.

Hence, in order to prove the claim, we must show that Eq. equation 11 holds. To this end, denote with $p_\theta(x, e)$ the deterministic decoder. Define the random variables

$$\begin{aligned}
\hat{X}_0 &:= p_\theta(x, e) \text{ and } E_0 := e, \text{ with } e \sim \mathcal{N}(e; 0, \sigma^2 I), \\
\hat{X}_\delta &:= p_\theta(x, e) \text{ and } E_\delta := e, \text{ with } e \sim \mathcal{N}(e; \delta, \sigma^2 I).
\end{aligned}$$

Note that in order to show that Eq. equation 11 holds, we need to show that

$$\begin{aligned}
\mathbb{P}((g \circ p_\theta)(x, e + \delta) = y_A) &\geq \mathbb{P}(g(\hat{X}_\delta) = y_A) \\
&> \mathbb{P}(g(\hat{X}_0) = y_A) = \mathbb{P}((g \circ p_\theta)(x, e + \delta) = y_B), \quad (12)
\end{aligned}$$

for each class $y_B \neq y_A$, and for each δ such that $\|\delta\| \leq R$. Without loss of generality, fix a class $y_B \neq y_A$.

For simplicity, for a fixed x we define $p_x^{-1}(\hat{x}) := e$ with e such that $p_\theta(x, e) = \hat{x}$. Define the half spaces

$$S_{\leq} := \{\delta^T p_x^{-1}(\hat{x}) \leq \sigma \|\delta\| \Phi^{-1}(\underline{p}_A)\},$$

$$S_{\geq} := \{\delta^T p_x^{-1}(\hat{x}) \geq \sigma \|\delta\| \Phi^{-1}(1 - \overline{p}_B)\}.$$

From Lemma B.3 it holds $\mathbb{P}(\hat{X}_0 \in S_{\leq}) = \underline{p}_A$, from which it follows that $\mathbb{P}(\hat{X}_0 \in S_{\leq}) = \underline{p}_A \leq p_A = \mathbb{P}(g(\hat{X}_0) = y_A)$. By combining this inequality with Lemma B.2 it holds

$$\mathbb{P}(\hat{X}_0 \in S_{\leq}) \leq \mathbb{P}(g(\hat{X}_0) = y_A) \Rightarrow \mathbb{P}(\hat{X}_\delta \in S_{\leq}) \leq \mathbb{P}(g(\hat{X}_\delta) = y_A) \quad (13)$$

Similarly, from Lemma B.3 it holds $\mathbb{P}(\hat{X}_0 \in S_{\geq}) = \overline{p}_B$, from which it follows that $\mathbb{P}(\hat{X}_0 \in S_{\geq}) = \overline{p}_B \geq p_B = \mathbb{P}(g(\hat{X}_0) = y_B)$. Again, by combining this inequality with Lemma B.2 it holds

$$\mathbb{P}(\hat{X}_0 \in S_{\geq}) \geq \mathbb{P}(g(\hat{X}_0) = y_A) \Rightarrow \mathbb{P}(\hat{X}_\delta \in S_{\geq}) \geq \mathbb{P}(g(\hat{X}_\delta) = y_A). \quad (14)$$

By Lemma B.5, since $\|\delta\| \leq R$ it holds $\mathbb{P}(\hat{X}_\delta \in S_{\leq}) \geq \mathbb{P}(\hat{X}_\delta \in S_{\geq})$. Then, it holds

$$\mathbb{P}(g(\hat{X}_\delta) = y_A) \geq \mathbb{P}(\hat{X}_\delta \in S_{\leq}) \geq \mathbb{P}(\hat{X}_\delta \in S_{\geq}) \geq \mathbb{P}(g(\hat{X}_\delta) = y_A),$$

where the first inequality follows from Eq. equation 13, the second one follows from Lemma B.5, and the last inequality follows from Eq. equation 14. Hence, Eq. equation 12 follows and so does the claim. \square

C PROOF OF THEOREM 4.4 AND LEMMA 4.5

Proof of Theorem 4.4. In order to prove the theorem, we give an explicit construction of the function p . Let L be a lower-triangular matrix that gives the Cholesky decomposition of Σ , i.e., $\Sigma = LL^T$, where L^T is the transpose of L . Then, we define our function p as

$$p(x) := \frac{1}{\sigma} Lx + \mu.$$

Note that the function p is deterministic and invertible, since the matrix L is positive-definite by definition of the Cholesky decomposition. Its inverse is given by the formula $p^{-1}(y) = \sigma L^{-1}(y - \mu)$. With this function, we can show that claims 1-2 hold.

(Proof of claim 1) The PDF of a multivariate Gaussian distribution $\mathcal{N}(x; 0, \sigma^2 I)$ is given by:

$$\mathbb{P}_X(x) = \frac{1}{\sqrt{(2\pi)^d \det(\sigma^2 I)}} \exp\left(-\frac{1}{2\sigma^2} \|x\|_2^2\right).$$

Similarly, the PDF of $\mathcal{N}(y; \mu, \Sigma)$ is:

$$\mathbb{P}_Y(y) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{1}{2}(y - \mu)^T \Sigma^{-1}(y - \mu)\right).$$

In order to prove the claim, we want to express $\mathbb{P}_Y(y)$ in terms of $\mathbb{P}_X(p^{-1}(y))$. To this end, by substituting $p^{-1}(y)$ into $\mathbb{P}_X(x)$ we get

$$\begin{aligned} \mathbb{P}_X(p^{-1}(y)) &= \frac{1}{\sqrt{(2\pi)^d \det(\sigma^2 I)}} \exp\left(-\frac{1}{2\sigma^2} \|p^{-1}(y)\|_2^2\right) \\ &= \frac{1}{\sqrt{(2\pi)^d \det(\sigma^2 I)}} \exp\left(-\frac{1}{2\sigma^2} \|\sigma L^{-1}(y - \mu)\|_2^2\right) \\ &= \frac{1}{\sqrt{(2\pi)^d \det(\sigma^2 I)}} \exp\left(-\frac{1}{2\sigma^2} (\sigma L^{-1}(y - \mu))^T (\sigma L^{-1}(y - \mu))\right) \\ &= \frac{1}{\sqrt{(2\pi)^d \det(\sigma^2 I)}} \exp\left(-\frac{1}{2}(y - \mu)^T (L^{-1})^T L^{-1}(y - \mu)\right) \\ &= \frac{1}{\sqrt{(2\pi)^d \det(\sigma^2 I)}} \exp\left(-\frac{1}{2}(y - \mu)^T (L^T L)^{-1}(y - \mu)\right) \\ &= \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{1}{2}(y - \mu)^T \Sigma^{-1}(y - \mu)\right) = \mathbb{P}_Y(y), \end{aligned}$$

where we have used that $\det(\sigma^2 I) = \det(\Sigma)$. Hence, claim 1 holds.

(Proof of claim 2) This claim follows directly from claim 1. In fact, since $\mathbb{P}_X(p^{-1}(y)) = \mathbb{P}_Y(y)$ and $\det(\sigma^2 I) = \det(\Sigma)$, we have that

$$-\frac{1}{2\sigma^2} \|p^{-1}(y)\|_2^2 = -\frac{1}{2}(y - \mu)^T \Sigma^{-1}(y - \mu).$$

The claim follows by standard calculations. \square

Proof. Proof of Lemma 4.5 First, note that it holds:

$$\det\left(\frac{\sigma^2}{\sqrt[d]{\det(\Sigma)}}\Sigma\right) = \frac{(\sigma^2)^d}{\det(\Sigma)} \det(\Sigma) = \det(\sigma^2 I).$$

Hence, by Thm. 4.4, there exists a deterministic invertible function p such that

$$\mathcal{N}\left(0, \frac{\sigma^2}{\sqrt[d]{\det(\Sigma)}}\Sigma\right) = p^\# \mathcal{N}(0, \sigma^2 I)$$

Hence, the claim follows by applying Thm. 4.3, and noting that it holds

$$\sqrt{(\hat{x} - x)^T \Sigma^{-1}(\hat{x} - x)} = \frac{1}{\sqrt[d]{\det(\Sigma)}} \|p^{-1}(y)\|_2.$$

\square

D WHITENING TRANSFORMATION

Whitening transformation is a linear transformation applied to a dataset to make the transformed variables uncorrelated and to standardize their variances (Kessy et al., 2018). Given a data matrix $X \in \mathbb{R}^{n \times d}$, where n is the number of samples and d is the number of features, whitening aims to transform X into a new matrix X_{whitened} such that the covariance matrix of X_{whitened} is the identity matrix I . The transformation is typically defined as $X_{\text{whitened}} = XW$, where $W \in \mathbb{R}^{d \times d}$ is the whitening matrix. The matrix W is derived from the covariance matrix Σ of the original data X , such that

$$\Sigma = \frac{1}{n} X^T X.$$

The primary goal is to find a matrix W that satisfies $W^T \Sigma W = I$. There are several methods to derive the whitening matrix W , each with its own practical advantages. In the following sub-sections we report on common methods to derive the matrix W . We refer the reader to, e.g., (Kessy et al., 2018) for a more comprehensive overview of these methods.

D.1 ZCA WHITENING (ZERO-PHASE COMPONENT ANALYSIS).

ZCA whitening aims to find a transformation that minimally alters the original data while achieving whitening. The matrix W is derived as $W_{\text{ZCA}} = Q\Lambda^{-1/2}Q^T$, where Q is the matrix of eigenvectors of the covariance matrix Σ , and Λ is a diagonal matrix containing the corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_d$ of Σ .

Practical Advantages.

- *Minimal Distortion:* ZCA whitening minimally distorts the original data in the least-squares sense, preserving the overall structure and appearance of the data.
- *Interpretability:* Since ZCA whitening preserves the spatial structure of the original data (e.g., images), it is often more interpretable in visual tasks.

Table 3: Base classifiers and execution time of RANDSMOOTH on CIFAR-10. In this table, **Params.** denotes the number of parameters (in millions), and **Time** denotes the average execution time required for a model to certify a datapoint, using the DCERT algorithm as in App. A, with parameters $\alpha = 0.001$, $n_0 = 100$ Monte Carlo samples for selection and $n = 100000$ samples for estimation (in seconds).

Model	Params. (m)	Time (s)	Model	Params. (m)	Time (s)
resnet20	0.27	2.1 ± 0.07	mobilenetv2_x0_5	0.7	5.61 ± 0.1
resnet32	0.47	3.46 ± 0.08	mobilenetv2_x1_4	4.33	15.36 ± 0.08
resnet44	0.66	4.53 ± 0.08	shufflenetv2_x1_0	1.26	4.56 ± 0.08
resnet56	0.86	5.73 ± 0.05	shufflenetv2_x0_5	0.35	2.53 ± 0.21
vgg13_bn	9.94	5.51 ± 0.1	shufflenetv2_x2_0	5.37	11.4 ± 0.07
vgg16_bn	15.25	7.05 ± 0.08	repvgg_a0	7.84	16.7 ± 0.07
vgg19_bn	20.57	9.04 ± 0.09	repvgg_a1	12.82	24.05 ± 0.18
mobilenetv2_x1_0	2.24	11.17 ± 0.07	repvgg_a2	26.82	25.17 ± 0.33

D.1.1 PCA WHITENING (PRINCIPAL COMPONENT ANALYSIS)

PCA Whitening involves projecting the data onto the principal components and scaling each component by the inverse square root of its corresponding eigenvalue. This method transforms the data so that it is uncorrelated and each principal component has unit variance. The whitening matrix W for PCA whitening is derived from the eigendecomposition of the covariance matrix Σ . The whitening matrix W in PCA whitening is computed as $W_{\text{PCA}} = \Lambda^{-1/2}Q^T$, where $\Lambda^{-1/2}$ is a diagonal matrix whose elements are the inverse square roots of the eigenvalues of Σ , and Q is the matrix whose columns are the eigenvectors of the covariance matrix Σ .

Practical Advantages.

- *Dimensionality Reduction*: PCA whitening naturally combines whitening with dimensionality reduction, as it allows discarding components corresponding to small eigenvalues, which may represent noise.
- *Variance Preservation*: It maximizes variance along the orthogonal axes, which is useful in scenarios where retaining variance in the principal directions is important.

D.1.2 CHOLESKY WHITENING

Cholesky whitening uses the Cholesky decomposition of the inverse covariance matrix. The whitening matrix W is computed as $W_{\text{Cholesky}} = L^{-1}$, where L is the lower triangular matrix from the Cholesky decomposition of the covariance matrix Σ (i.e., $\Sigma = LL^T$).

Practical Advantages.

- *Computational Efficiency*: Cholesky whitening is computationally efficient for large-scale datasets because it leverages a triangular decomposition, which is faster to compute than full eigenvalue decomposition.
- *Numerical Stability*: This method is numerically stable, especially when Σ is well-conditioned.

E ADDITIONAL EXPERIMENTS

E.1 EXPERIMENTS ON EXECUTION TIME

We provide tables on the execution time of RANDSMOOTH on CIFAR-10 (Table 3) and IMAGENET (Table 4). Overall, we observe that the execution time of RANDSMOOTH is similar to the execution time of DSMOOTH. Furthermore, we provide tables on the execution time of LSMOOTH on CIFAR-10 (Table 5) and IMAGENET (Table 6). Again, we observe that the execution time of RANDSMOOTH is similar to the execution time of DSMOOTH.

E.2 EXPERIMENTS ON CERTIFIED ACCURACY

We provide additional experiments on the certified accuracy of the base classifiers in Table 1 on CIFAR-10, for various parameters choices.

Table 4: Base classifiers and execution time of RANDSMOOTH on IMAGENET. In this table, **Params.** denotes the number of parameters (in millions), and **Time** denotes the average execution time required for a model to certify a datapoint, using the DCERT algorithm as in App. A , with parameters $\alpha = 0.001$, $n_0 = 100$ Monte Carlo samples for selection and $n = 1000$ samples for estimation. (in seconds).

Model	Params. (m)	Time (s)	Model	Params. (m)	Time (s)
resnet50	25.56	26.3 ± 0.44	wide_resnet50_2	68.88	40.02 ± 0.64
resnet152	60.19	169.84 ± 8.01	wide_resnet101_2	126.89	85.45 ± 11.73

Table 5: Base classifiers and execution time of LSMOOTH on CIFAR-10. In this table, **Params.** denotes the number of parameters (in millions), and **Time** denotes the average execution time required for a model to certify a datapoint, using the DCERT algorithm as in App. A , with parameters $\alpha = 0.001$, $n_0 = 100$ Monte Carlo samples for selection and $n = 100000$ samples for estimation (in seconds).

Model	Params. (m)	Time (s)	Model	Params. (m)	Time (s)
resnet20	0.27	51.42 ± 2.2	mobilenetv2_x0_5	0.7	56.52 ± 3.21
resnet32	0.47	52.92 ± 5.05	mobilenetv2_x1_4	4.33	64.82 ± 6.11
resnet44	0.66	45.05 ± 24.42	shufflenetv2_x1_0	1.26	25.89 ± 3.72
resnet56	0.86	32.43 ± 2.12	shufflenetv2_x0_5	0.35	9.4 ± 3.8
vgg13_bn	9.94	8.53 ± 0.97	shufflenetv2_x2_0	5.37	28.38 ± 5.32
vgg16_bn	15.25	24.8 ± 4.57	repvgg_a0	7.84	29.44 ± 1.63
vgg19_bn	20.57	11.57 ± 1.32	repvgg_a1	12.82	26.97 ± 2.15
mobilenetv2_x1_0	2.24	15.88 ± 3.49	repvgg_a2	26.82	39.19 ± 2.25

Effect of the number n of samples for estimation. We perform a set of experiments to determine the effects of the number n of samples for estimation on the performance of DSMOOTH, which is quantified using the approximate certified accuracy. To this end, we run our algorithm with parameters $\alpha = 0.001$, $n_0 = 100$ Monte Carlo samples for selection, and $n = 100, 1000, 10000, 100000, 1000000$ samples for estimation. The results are displayed in Fig. 4. We observe, that the number of samples affects the certified accuracy. This results highlight a well-known feature of smoothing algorithms, such as DSMOOTH, namely that these algorithms require a large number of samples to achieve good certified accuracy.

Effect of the parameter α . We conduct a series of experiments to evaluate how the parameter α influences the performance of DSMOOTH, measured by the approximate certified accuracy. Our testing involves running the algorithm with $n_0 = 100$ Monte Carlo samples for selection, $n = 100000$ samples for estimation, and varying α values of 0.1, 0.01, 0.001, 0.0001, and 0.00001. The outcomes are shown in Fig. 5. We find that while the number of samples impacts the certified accuracy, the parameter α does not notably affect the model’s performance.

F EXAMPLES OF NOISY IMAGES

We show examples of CIFAR-10 (Fig. 6) and IMAGENET (Fig. 7) images corrupted with the smoothing distribution as in equation 3.

Table 6: Base classifiers and execution time of LSMOOTH on IMAGENET. In this table, **Params.** denotes the number of parameters (in millions), and **Time** denotes the average execution time required for a model to certify a datapoint, using the DCERT algorithm as in App. A, with parameters $\alpha = 0.001$, $n_0 = 100$ Monte Carlo samples for selection and $n = 1000$ samples for estimation. (in seconds).

Model	Params. (m)	Time (s)	Model	Params. (m)	Time (s)
resnet50	25.56	27.83 ± 0.18	wide_resnet50_2	68.88	41.0 ± 0.17
resnet152	60.19	89.53 ± 2.62	wide_resnet101_2	126.89	68.08 ± 1.67

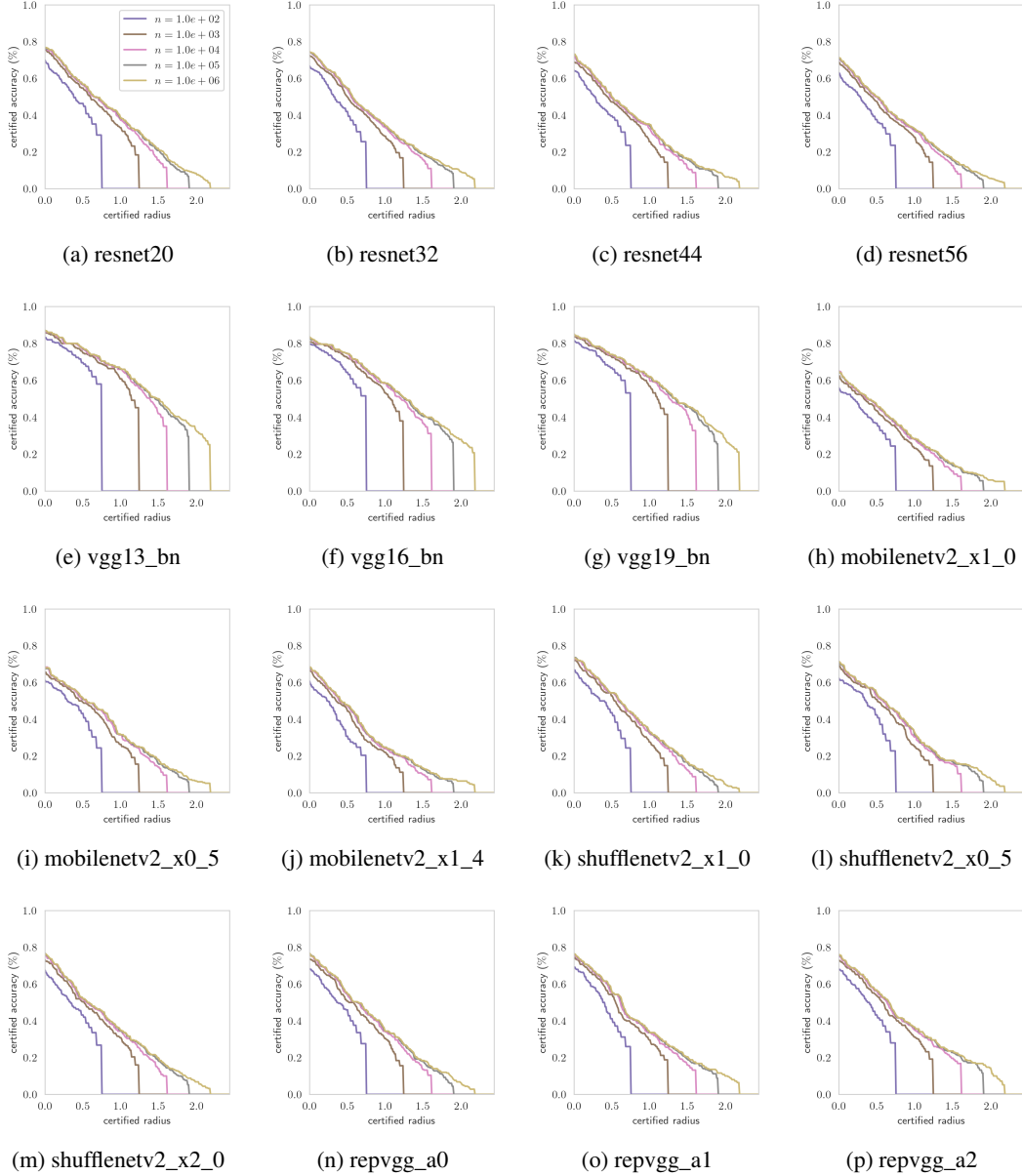


Figure 4: Approximate certified accuracy of DCERT on various base models as in Table 1. We use parameters $\alpha = 0.001$ and $n_0 = 100$ Monte Carlo samples for selection. We show the certified accuracy for various choices of the number n of samples for estimation.

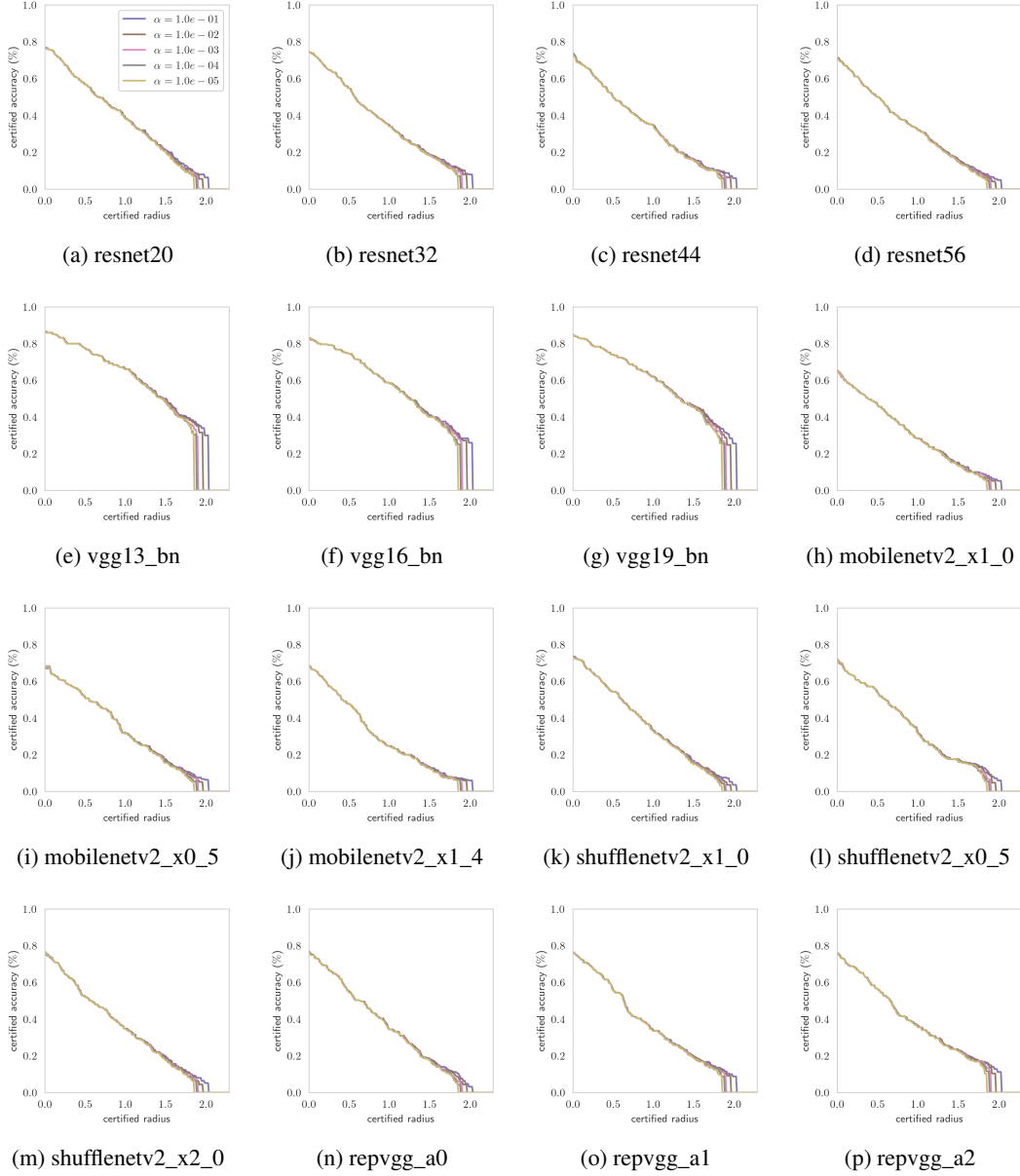


Figure 5: Approximate certified accuracy of DCERT on various base models as in Table 1. We use parameters $n_0 = 100$ Monte Carlo samples for selection, and $n = 100000$ samples for estimation. We show the certified accuracy for various choices of the parameter α .

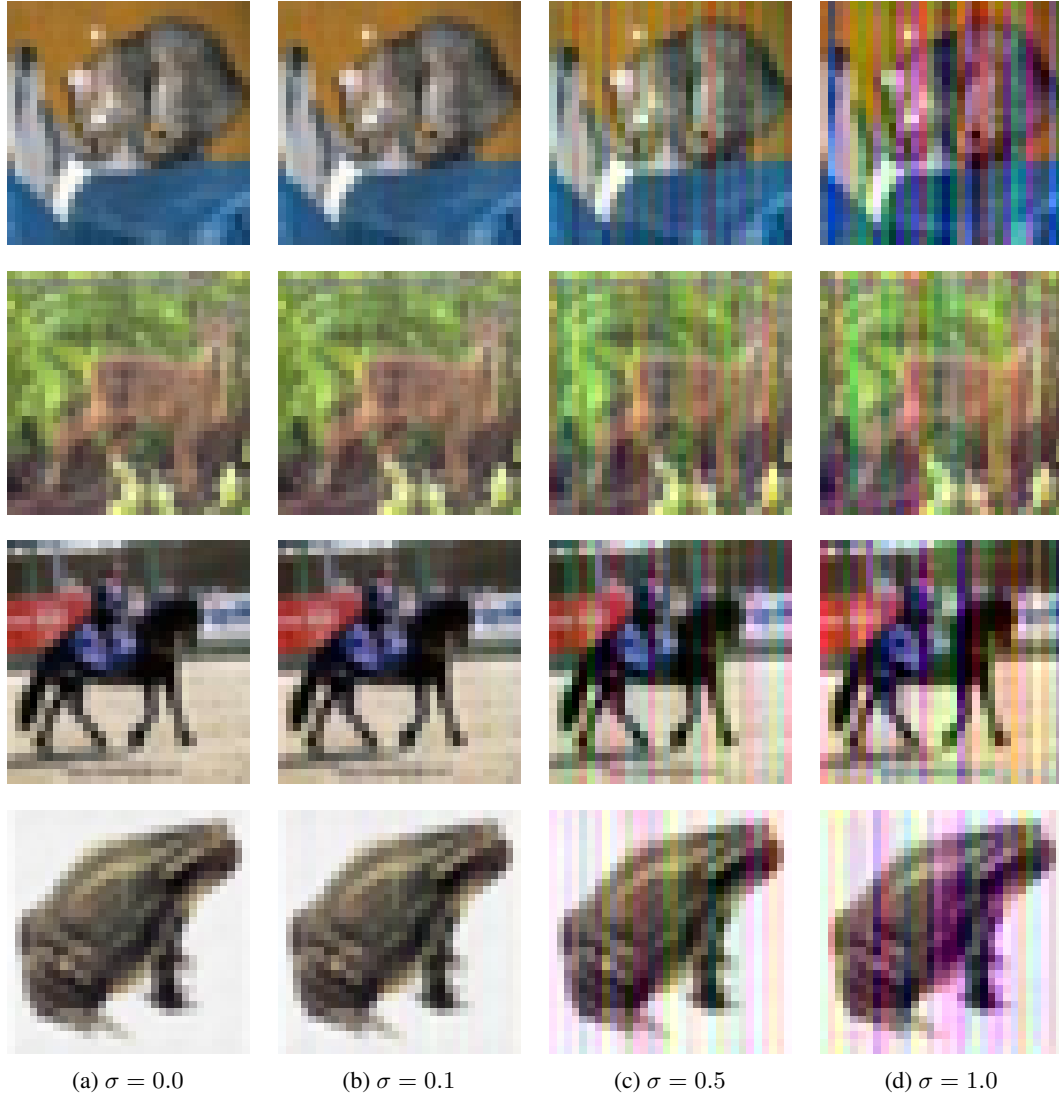


Figure 6: Examples of CIFAR-10 images corrupted with the smoothing distribution as in equation 3 for a SQUARE + FGSM attack as described in Sec. 3.1. Here, we use resnet32 as a base classifier.

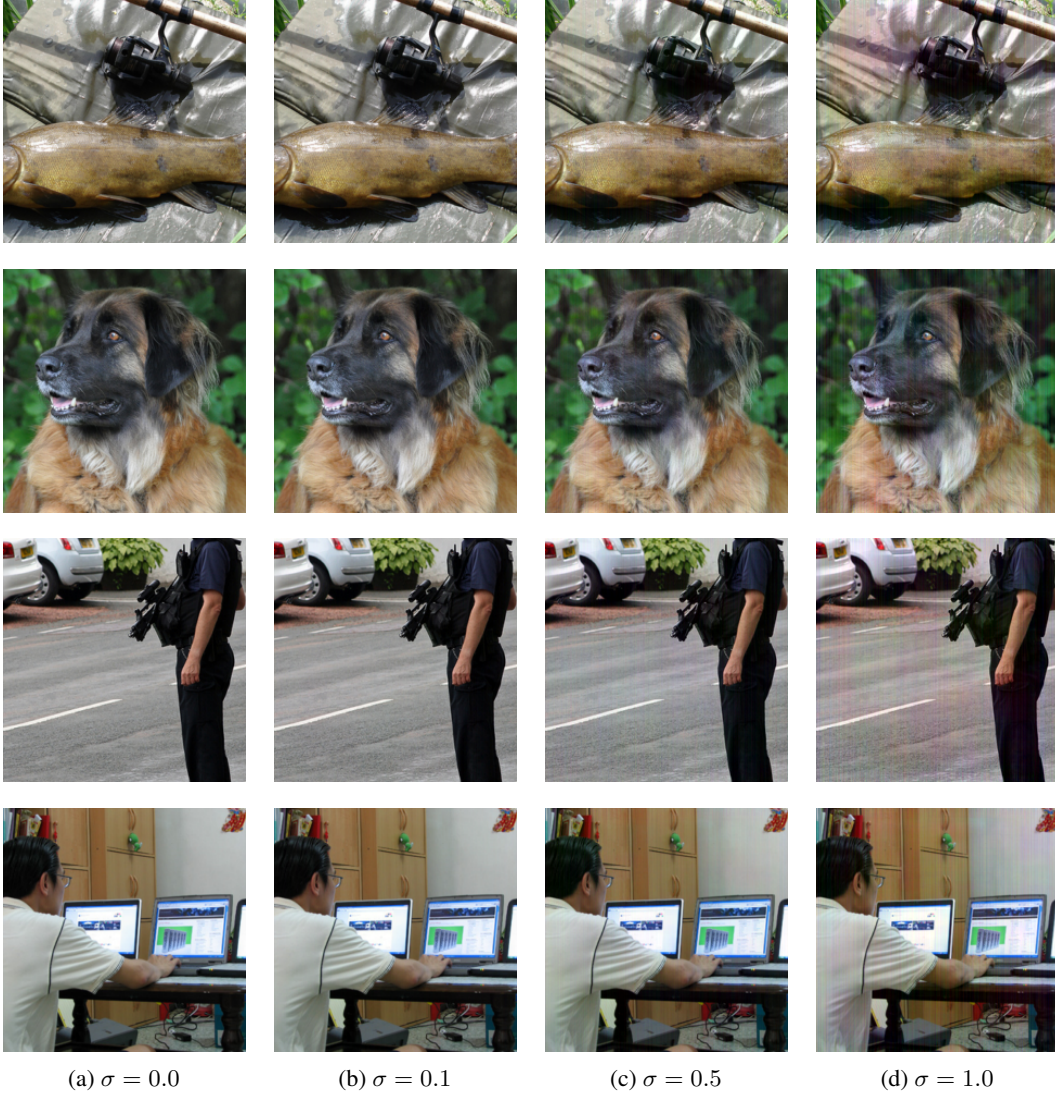


Figure 7: Examples of IMAGENET images corrupted with the smoothing distribution as in equation 3 for a SQUARE + FGSM attack as described in Sec. 3.1. Here, we use resnet50 as a base classifier.