

# Input-Specific Robustness Certification for Randomized Smoothing

Ruoxin Chen<sup>1</sup>, Jie Li<sup>1</sup>\*, Junchi Yan<sup>1</sup>, Ping Li<sup>2</sup>, Bin Sheng<sup>1</sup>

<sup>1</sup> Shanghai Jiao Tong University

<sup>2</sup> The Hong Kong Polytechnic University

## Abstract

Although randomized smoothing has demonstrated high certified robustness and superior scalability to other certified defenses, the high computational overhead of the robustness certification bottlenecks the practical applicability, as it depends heavily on the large sample approximation for estimating the confidence interval. In existing works, the sample size for the confidence interval is universally set and agnostic to the input for prediction. This Input-Agnostic Sampling (IAS) scheme may yield a poor Average Certified Radius (ACR)-runtime trade-off which calls for improvement. In this paper, we propose Input-Specific Sampling (ISS) acceleration to achieve the cost-effectiveness for robustness certification, in an adaptive way of reducing the sampling size based on the input characteristic. Furthermore, our method universally controls the certified radius decline from the ISS sample size reduction. The empirical results on CIFAR-10 and ImageNet show that ISS can speed up the certification by more than three times at a limited cost of 0.05 certified radius. Meanwhile, ISS surpasses IAS on the average certified radius across the extensive hyper-parameter settings. Specifically, ISS achieves ACR=0.958 on ImageNet ( $\sigma = 1.0$ ) in 250 minutes, compared to ACR=0.917 by IAS under the same condition. We release our code in <https://github.com/roy-ch/Input-Specific-Certification>.

## 1 Introduction

Neural networks are known susceptible to adversarial attacks (Szegedy et al. 2014; Goodfellow, Shlens, and Szegedy 2015). A line of empirical defenses (Buckman et al. 2018; Song et al. 2018) have been proposed to defend adversarial attacks, but are often broken by the newly devised stronger attacks (Athalye, Carlini, and Wagner 2018). Existing certified defenses (Wong et al. 2018; Raghu et al. 2018; Liang et al. 2018; Cohen, Rosenfeld, and Kolter 2019) provide the theoretical guarantees for their robustness. In particular, *Randomized smoothing* (Cohen, Rosenfeld, and Kolter 2019) is one of the few certified defenses that can scale to ImageNet-scale classification task, showing its great potential for wide application. Moreover, randomized smoothing has shown high robustness

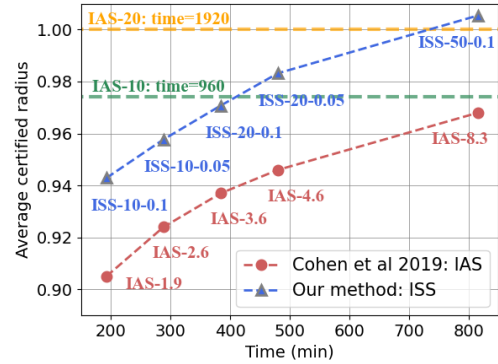


Figure 1: ISS achieves a better ACR-time trade-off than IAS. ISS -  $c_1$  -  $c_2$  denotes that ISS accelerates IAS -  $c_1$  at the controllable certified radius decline ( $\leq c_2$ ). IAS -  $c_1$  denotes that IAS accelerates IAS - 10, IAS - 20 by reducing the sample size to  $c_1 \times 10,000$ . ISS always surpasses IAS on ACR in the same time cost. The results are evaluated on the ImageNet ( $\sigma = 1.0$ ) model trained by (Jeong and Shin 2020).

against various types of adversarial attacks, including norm-constrained perturbations (e.g.  $\ell_0$ ,  $\ell_2$ ,  $\ell_\infty$  norms) and image transformations (e.g. rotations and image shift).

Despite these advances, randomized smoothing suffers the costly robustness certification. Specifically, computing a certified radius close to the exact value needs a relatively tight lower bound of the top-1 label probability, which requires running forward passes on a large number of samples (Salman et al. 2019; Cohen, Rosenfeld, and Kolter 2019; Zhai et al. 2020; Jeong and Shin 2020; Jia et al. 2020). Such expensive overheads make them less applicable to the real-world scenarios. Some works (Jia et al. 2020; Feng et al. 2020) proposed to leverage the runner-up label probability in the certification, but their performances may suffer from the inevitable loss in the simultaneous confidence intervals. Traditionally, the robustness certification is accelerated by reducing the sample size used for estimating the lower bound (Cohen, Rosenfeld, and Kolter 2019; Jia et al. 2020), but the vanilla sample size reduction will lead to a poor ACR-runtime trade-off. It is critical to develop a cost-effective certification method.

In this paper, we propose Input-Specific Sampling (ISS) to speed up the certification for randomized smoothing, with-

\*Jie Li is the corresponding author, who is with the Department of Computer Science and Engineering and MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University Shanghai, China. Email: lijiecs@sjtu.edu.cn  
Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

out hurting too much on the certification performance. The idea behind ISS is to *minimize the sample size for the given input at the bounded cost of the certified radius decline*, instead of directly applying the same sample size to all inputs. The idea is realized by precomputing a mapping from the input characteristics to the sample size. Consequently, ISS can accelerate the certification at a controllable cost. Empirical results validate that ISS consistently outperforms IAS (Cohen, Rosenfeld, and Kolter 2019) on ACR. As shown in Fig. 1, ISS – 10 – 0.05 (ACR=0.958) accelerates the standard certification IAS – 10, shortening the certification time 962 → 250 mins at the controllable decline ( $\leq 0.05$ ). Furthermore, ISS is compatible with all the randomized smoothing works that need confidence interval, since ISS has no additional constraint on the base classifier or the smoothing scheme.

Our contributions can be summarized as follows;

1. We propose Input-Specific Sampling (ISS) to adaptively reduce the sample size for each input. The proposed input-specific sampling, for the first time to our best knowledge, can significantly reduce the cost for accelerating the robustness certification of randomized smoothing.
2. ISS can universally control the difference between the certified radii before and after the acceleration. In particular, the sample size computed by ISS is theoretically tight for bounding the radius decline.
3. The results on CIFAR-10 and ImageNet demonstrate that: 1) ISS significantly accelerates the certification at a controllable decline in the certified radii. 2) ISS consistently achieves a higher average certified radius when compared to the mainstream acceleration IAS.

## 2 Related Works

**Certified defenses.** Neural networks are vulnerable to adversarial attacks (Athalye, Carlini, and Wagner 2018; Eykholt et al. 2018; Kurakin, Goodfellow, and Bengio 2017; Eykholt et al. 2018; Jia and Gong 2018). Compared to empirical defenses (Goodfellow, Shlens, and Szegedy 2015; Svoboda et al. 2019; Buckman et al. 2018; Ma et al. 2018; Guo et al. 2018; Dhillon et al. 2018; Xie et al. 2017; Song et al. 2018), certified defenses can provide provable robustness guarantees for their predictions. Recently, a line of certified defenses have been proposed, including dual network (Wong et al. 2018), convex polytope (Wong and Kolter 2018), CROWN-IBP (Zhang et al. 2019), Lipschitz bounding (Cisse et al. 2017). However, those certified defenses suffer from either the low scalability or the hard constraints on the neural network architecture.

**Randomized smoothing.** In the seminal work (Cohen, Rosenfeld, and Kolter 2019), the authors for the first time propose randomized smoothing to defend the  $\ell_2$ -norm perturbations, which significantly outperforms other certified defenses. Recently, series of works further extend randomized smoothing to defend various attacks, including  $\ell_0, \ell_1, \ell_2, \ell_\infty$ -norm perturbations and geometric transformations. For instance, (Levine and Feizi 2020) introduce the random ablation against  $\ell_0$ -norm adversarial attacks. (Yang et al. 2020) propose Wulff Crystal uniform distribution against  $\ell_1$ -norm perturbations. (Awasthi et al. 2020) introduce  $\infty \rightarrow 2$  matrix operator for Gaussian smoothing to defend  $\ell_\infty$ -norm

perturbations. (Fischer, Baader, and Vechev 2020; Li et al. 2020) exploit randomized smoothing to defend adversarial translations. Remarkably, almost all the randomized smoothing works (Salman et al. 2019; Cohen, Rosenfeld, and Kolter 2019; Zhai et al. 2020; Jeong and Shin 2020; Yang et al. 2020; Jia et al. 2020) have achieved superior certified robustness to other certified defenses in their respective fields.

**Robustness certification in randomized smoothing.** Despite its sound performance, the certification of randomized smoothing is seriously costly. Unfortunately, accelerating the certification is a fairly under-explored field. The mainstream acceleration method (Jia et al. 2020; Feng et al. 2020), which we call IAS, is to apply a smaller sample size for certifying the radius. However, IAS accelerates the certification at a seriously sacrifice ACR and the certified radii of specific inputs. Therefore, it calls for approaches to achieve a better time-cost trade-off, which is the main purpose of this paper.

## 3 Preliminaries

**Randomized smoothing** The basic idea of randomized smoothing (Cohen, Rosenfeld, and Kolter 2019) is to generate a smoothed version of the base classifier  $f$ . Given an arbitrary base classifier  $f(x) : \mathbb{R}^d \rightarrow \mathcal{Y}$  where  $\mathcal{Y} = \{1, \dots, k\}$  is the output space, the smoothed classifier  $g(\cdot)$  is defined as:

$$g(x) := \arg \max_{c \in \mathcal{Y}} Pr[f(x') = c], \quad x' \sim \mathcal{N}(x, \sigma^2 I^d) \quad (1)$$

$g(x)$  returns the most likely predicted label of  $f(\cdot)$  when input the data with Gaussian augmentation  $x' \sim \mathcal{N}(x, \sigma^2 I^d)$ . The tight lower bound of  $\ell_2$ -norm certified radius (Cohen, Rosenfeld, and Kolter 2019) for the prediction  $c_A = g(x)$  is:

$$\sigma \Phi^{-1}(p_A) \quad \text{where } p_A := Pr[f(x') = c_A], \quad x' \sim \mathcal{N}(x, \sigma^2 I^d) \quad (2)$$

where  $\Phi^{-1}$  is the inverse of the standard Gaussian CDF. We emphasize that computing the deterministic value of  $g(x)$  is impossible because  $g(\cdot)$  is built over the random distribution  $\mathcal{N}(x, \sigma^2 I^d)$ . Therefore, we use Clopper-Pearson method (Clopper and Pearson 1934) to guarantee  $Pr[f(x') = c_A] > Pr[f(x') = c], \forall c \neq c_A$  with the confidence level  $1 - \alpha$ , and then we have  $g(x) = c_A$  with the confidence level  $1 - \alpha$ .

**Robustness certification** In practice, the main challenge in computing the radius  $\sigma \Phi^{-1}(p_A)$  is that  $p_A$  is inaccessible because iterating all possible  $f(x') : x' \in \mathbb{R}^d$  is impossible. Therefore, we estimate  $\underline{p}_A$ , the standard one-sided Clopper-Pearson confidence lower bound of  $p_A$  instead of  $p_A$  and certify a lower bound  $\sigma \Phi^{-1}(\underline{p}_A)$ . Estimating a tight  $\underline{p}_A$  needs a large size of samples for  $f(x') : x' \sim \mathcal{N}(x, \sigma^2 I^d)$ . Generally, the estimated  $\underline{p}_A$  increases with the sample size<sup>1</sup>.

<sup>1</sup>The seminal work (Cohen, Rosenfeld, and Kolter 2019) derives the certified radius:  $\frac{\sigma}{2} [\Phi^{-1}(p_A) - \Phi^{-1}(p_B)]$  where  $p_B$  is the runner-up label probability. Currently, most works (Cohen, Rosenfeld, and Kolter 2019; Zhai et al. 2020; Jeong and Shin 2020; Jia et al. 2020) compute the certified radius by Eq. (2), which substitutes  $p_B$  with  $1 - p_A$ , to avoid doing interval estimation twice.

## Standard certification and vanilla acceleration: IAS

The standard certification algorithm (Cohen, Rosenfeld, and Kolter 2019) can be summarized in two steps:

1. **Sampling:** Given the input  $x$ , sample  $k$  (e.g.  $k = 100,000$ ) iid samples  $\{x'_i : i = 1, \dots, k\} \sim \mathcal{N}(x, \sigma^2 I^d)$  and run  $k$  times forward passes  $\{f(x'_i) : i = 1, \dots, k\}$ .
2. **Interval estimation:** Count  $k_A = \sum_{i=1}^k \mathbb{I}\{f(x'_i) = c_A\}$  ( $\mathbb{I}$  denotes the indicator function) where  $c_A$  is the label with top-1 label counts. Compute the confidence lower bound  $p_A$  with the confidence level  $1 - \alpha$ . Return the certified radius  $\sigma \Phi^{-1}(p_A)$ .

The high computation is mainly due to the  $k$  times forward passes in **Sampling**. The certification is accelerated by the vanilla sample size reduction, which we call input-agnostic sample size reduction (IAS). This acceleration is at the cost of unpredictable radius declines, which yields a poor ACR-runtime trade-off since it reduces the sample size equally for each input, without considering the input characteristics.

## 4 Methodology

We first introduce the notions of Absolute Decline and Relative Decline. Then we propose Input-Specific Sampling (ISS), which aims to use the minimum sample size with the constraint that the radius decline is less than the given bound.

### 4.1 Overview and main idea

The key idea of ISS is to appropriately reduce the sample size for each input, instead of applying the same sample size to the certifications for all inputs. Since the sample size reduction will inevitably cause the decline in the certified radius, thus we aim to quantify the radius decline and bound the decline to be less than the pre-specified value. First we define the radius decline as follows:

**Definition 1 (Absolute Decline  $\text{AD}(k; \bar{k}, p)$ ).** Given the input  $x$  and the pre-specified desired sample size  $\bar{k}$  (e.g.  $\bar{k} = 100,000$ ), suppose we know  $p_A$  of  $x$ , Absolute Decline  $\text{AD}(k; \bar{k}, p)$  is the gap between the radius certified at the sample size  $\bar{k}$  and the radius certified at  $k : k \leq \bar{k}$ :

$$\text{AD}(k; \bar{k}, p_A) := \underbrace{\sigma \Phi^{-1}(p_1)}_{\text{Desired radius}} - \underbrace{\sigma \Phi^{-1}(p_2)}_{\text{Estimated radius}} \quad (3)$$

where  $p_1 = \mathbf{B}(\alpha; p_A \bar{k}, \bar{k} - p_A \bar{k} + 1)$ ,  
 $p_2 = \mathbf{B}(\alpha; p_A k, k - p_A k + 1)$

where  $\mathbf{B}(\alpha; k_A, k - k_A + 1)$  denotes the one-sided Clopper-Pearson lower bound (Clopper and Pearson 1934) with the confidence level  $1 - \alpha$ , which is equal to the  $\alpha$ th quantile from a Beta distribution with shape parameters  $k_A, k - k_A + 1$ .

**Definition 2 (Relative Decline  $\text{RD}(k; \bar{k}, p_A)$ ).** Similar to absolute decline, Relative Decline  $\text{RD}(k; \bar{k}, p_A)$  is

$$\text{RD}(k; \bar{k}, p_A) := \frac{\sigma \Phi^{-1}(p_1) - \sigma \Phi^{-1}(p_2)}{\sigma \Phi^{-1}(p_1)} \quad (4)$$

where  $p_1 = \mathbf{B}(\alpha; p_A \bar{k}, \bar{k} - p_A \bar{k} + 1)$ ,  
 $p_2 = \mathbf{B}(\alpha; p_A k, k - p_A k + 1)$

---

### Algorithm 1: Compute ISS mapping $\psi_{\text{ISS}}(\cdot)$

---

**Input:** The maximum decline  $\bar{U}$ , the decline type, the desired sample size  $\bar{k}$ , the noise level  $\sigma$ , the confidence level  $1 - \alpha$ , the length of the subinterval  $\delta$

**Output:** the ISS mapping  $\psi_{\text{ISS}}(\cdot)$

```

1: for  $N = 0, 1, 2, \dots, \frac{1}{\delta}$  do
2:    $p \leftarrow N \cdot \delta$ ;
3:   Compute  $\bar{r} = \sigma \cdot \Phi^{-1}(\mathbf{B}(\alpha; p\bar{k}, \bar{k}))$ ;
4:   Compute the minimum required certified radius:
     If the decline type is AD:  $\tilde{r} \leftarrow \bar{r} - U_{\text{AD}}$  or
     If the decline type is RD:  $\tilde{r} \leftarrow (1 - U_{\text{RD}})\bar{r}$ ;
5:   if  $\tilde{r} \leq 0$  then
6:      $\psi_{\text{ISS}}(p) \leftarrow 0$ ;
7:   else
8:      $\psi_{\text{ISS}}(p) \leftarrow \arg \min_k \sigma \cdot \Phi^{-1}(\mathbf{B}(\alpha; pk, k)) \geq \tilde{r}$ ;
9:   end if
10: end for
11: Return  $\psi_{\text{ISS}}(p) : p = \delta, 2\delta, \dots, 1$ ;

```

---

**Remark 1.** The absolute (or relative) decline is the expected gap between the radius certified at the sample size  $\bar{k}$  and  $k$  when fixing  $k_A/k \equiv p_A$  where  $k_A := \sum_{i=1}^k \mathbb{I}\{f(x'_i) = c_A\}$ . It connects the expected radius decline to the sample size when given  $p_A$ . In particular, when  $\bar{k} = \infty$ , the absolute (or relative) decline measures the gap between the optimal certified radius that randomized smoothing can provide and the radius certified at the sample size  $k$ .

**Formulate our key idea** Given the input  $x$  and the pre-specified upper bound of the decline  $U_{\text{AD}} \in \mathbb{R}^+$  (or  $U_{\text{RD}} \in \mathbb{R}^+$ ), our idea for AD (or RD) is formulated as follows:

1. find  $\min k$  with the constraint  $\text{AD}(k; \bar{k}, p) \leq U_{\text{AD}}$ .
2. find  $\min k$  with the constraint  $\text{RD}(k; \bar{k}, p) \leq U_{\text{RD}}$ .

In practice, solving the above two problems is non-trivial because  $p_A$  of  $x$  is inaccessible. Simply treating the estimated  $k_A/k$  as  $p_A$  is obviously unreasonable. We propose ISS, a practical solution to the above two problems.

### 4.2 Certification with input-specific sampling

Fig. 2 shows an overview. Given the input  $x$ , we first estimate a relatively loose two-sided Clopper-Pearson confidence interval  $p_A \in [p_{\text{low}}, p_{\text{up}}]$  by  $k_0$  samples where  $k_0 < \bar{k}$  is a relatively small sample size. Given  $\bar{k}, U_{\text{AD}}$  (or  $U_{\text{RD}}$ ), ISS assigns the sample size  $\hat{k}$  for certifying  $g(x)$  where  $\hat{k}$  is:

$$\hat{k} = \max(\psi(p_{\text{low}}), \psi(p_{\text{up}})) \quad (5)$$

For Absolute Decline :  $\psi(p) := \arg \min_k \text{AD}(k; \bar{k}, p) \leq U_{\text{AD}}$

For Relative Decline :  $\psi(p) := \arg \min_k \text{RD}(k; \bar{k}, p) \leq U_{\text{RD}}$

Formally, we present the following two propositions to theoretically prove that  $\hat{k}$  (AD) computed from Eq. (5) is optimal. Prop. 1 guarantees that the sample size  $\hat{k}$  computed from Eq. (5) must satisfy the constraint  $\text{AD}(\hat{k}; \bar{k}, p_A) \leq$

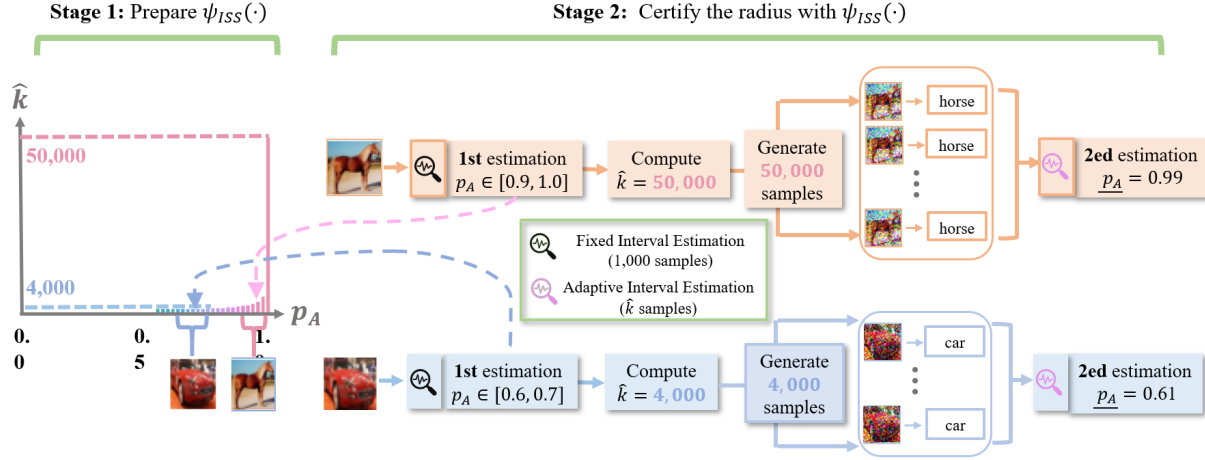


Figure 2: Overview of the robustness certification with ISS. In Stage 1, we compute  $\psi_{ISS}(\cdot)$ , a mapping from  $p_A$  to  $\hat{k}$ . In Stage 2, given the image  $x$ , we first loosely estimate the confidence interval for  $p_A$  and determine the sample size for  $x$  with  $\psi_{ISS}(\cdot)$ .

---

**Algorithm 2: Certification with input-specific sampling (ISS)**

---

**Input:** The input  $x$ , the base classifier  $f$ , the maximum sample size  $\bar{k}$ , the sample size  $k_0 : k_0 \leq \bar{k}$ , the confidence level  $\alpha$ , the ISS mapping  $\psi_{ISS}(\cdot)$

**Output:** Prediction  $\text{pred}$ , radius  $r$

- 1: Sample  $k_0$  noisy samples  $x'_1, \dots, x'_k \sim \mathcal{N}(x, \sigma^2 I^d)$ ;
  - 2: Compute the prediction:  
 $\text{pred} \leftarrow \arg \max_{y \in \mathcal{Y}} \sum_{i=1}^{k_0} \mathbb{I}\{f(x'_i) = y\}$ ;
  - 3: Count  $k_A^0 \leftarrow \max_{y \in \mathcal{Y}} \sum_{i=1}^{k_0} \mathbb{I}\{f(x'_i) = y\}$ ;
  - 4: Compute the two-sided confidence interval:  
 $p_{\text{low}} \leftarrow \mathbf{B}(\alpha/2; k_A^0, k_0 - k_A^0 + 1)$   
 $p_{\text{up}} \leftarrow \mathbf{B}(1 - \alpha/2; k_A^0 + 1, k_0 - k_A^0)$ ;
  - 5: Compute  $\hat{k} \leftarrow \max(\psi_{ISS}(p_{\text{low}}), \psi_{ISS}(p_{\text{up}}))$ ;
  - 6: Sample  $\max(\hat{k} - k_0, 0)$  noisy samples:  
 $x'_{k_0+1}, \dots, x'_k \sim \mathcal{N}(x, \sigma^2 I^d)$ ;
  - 7: Count  $k_A \leftarrow \max_{y \in \mathcal{Y}} \sum_{i=1}^{\hat{k}} \mathbb{I}\{f(x'_i) = y\}$ ;
  - 8: Compute the one-sided confidence lower bound:  
 $\underline{p}_A \leftarrow \mathbf{B}(\alpha; k_A, \hat{k} - k_A + 1)$ ;
  - 9: **if**  $\underline{p}_A < \frac{1}{2}$  **then**
  - 10:    $\text{pred} \leftarrow \text{ABSTAIN}$ ,  $r \leftarrow 0$ ;
  - 11: **else**
  - 12:   Compute the radius  $r \leftarrow \sigma \Phi^{-1}(\underline{p}_A)$ ;
  - 13: **end if**
  - 14: **Return**  $\text{pred}$  and  $r$ ;
- 

$U_{AD}$ . Prop. 2 guarantees that  $\hat{k}$  is the minimize sample size that can guarantee  $\text{AD}(\hat{k}; \bar{k}, p_A) \leq U_{AD}$ .

**Proposition 1. [Bounded absolute radius decline]** Suppose  $p_A \in [p_{\text{low}}, p_{\text{up}}]$  with  $1 - \alpha$  confidence level, then we guarantee that there is at least  $1 - \alpha$  probability that  $\hat{k}$  computed from Eq. (5) satisfies  $\text{AD}(\hat{k}; \bar{k}, p_A) \leq U_{AD}$ .

**Proposition 2. [Tightness for  $\hat{k}$ ]** Suppose  $p_A \in [p_{\text{low}}, p_{\text{up}}]$  and  $\hat{k}$  is computed from Eq. (5), then for an arbitrary sample size  $k : k < \hat{k}$ , there exists  $p_A \in [p_{\text{low}}, p_{\text{up}}]$  that breaks the

$$\text{constraint } \text{AD}(k; \bar{k}, p_A) \leq U_{AD}.$$

### 4.3 Implementation

In the practical algorithm of ISS, we substitute  $\psi(\cdot)$  in Eq. (5) with a piecewise constant function approximation  $\psi_{ISS}(\cdot)$ . The advantage of  $\psi_{ISS}(\cdot)$  over  $\psi(\cdot)$  is that we can compute  $\psi_{ISS}(p) : p \in [0, 1]$  previously before the certification to save the cost in computing  $\psi(p_{\text{low}}), \psi(p_{\text{up}})$  in Eq. (5) when certifying the radius for the testing data. Constructing  $\psi_{ISS}(\cdot)$  is feasible because that the value of  $\psi(p)$  only depends on  $p$  when fixing  $\bar{k}$ , regardless of the testing set or the base classifier architecture. Specifically,  $\psi_{ISS}(p)$  is

$$\psi_{ISS}(p) = \begin{cases} \psi(p) & p/\delta \in \mathbb{N} \\ \max(\psi(N_1\delta), \psi(N_2\delta)) & p/\delta \in (N_1, N_2) \end{cases} \quad (6)$$

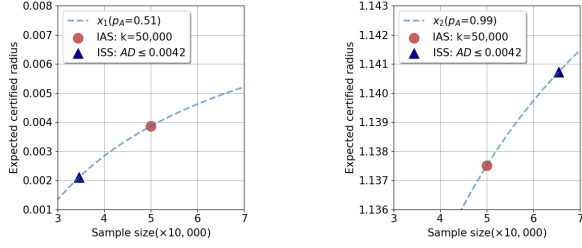
where  $N_1, N_2 \in \mathbb{N}$ . Obviously,  $\forall p \in [0, 1], \psi_{ISS}(p) \geq \psi(p)$ , thus Prop. 1 still holds for the substitution  $\psi_{ISS}(\cdot)$ . Prop. 2 holds for  $\psi_{ISS}(\cdot)$  when  $p_{\text{low}}/\delta \in \mathbb{N}, p_{\text{up}}/\delta \in \mathbb{N}$ .

The practical algorithm can summarized into two stages:

**Stage 1: prepare  $\psi_{ISS}(\cdot)$ .** Given  $\bar{k}$  and the decline upper bound  $U_{AD}$  (or  $U_{RD}$ ), compute  $\psi_{ISS}(p)$  by Eq. (5) and Eq. (6). The detailed algorithm is shown in Alg. 1. **Stage 2: cer-**

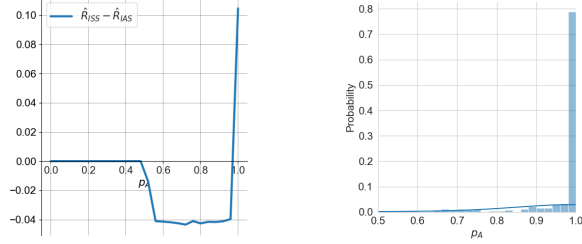
**tify the radius with  $\psi_{ISS}(\cdot)$ .** Given  $x$ , we first estimate a loose confidence interval  $p_A \in [p_{\text{low}}, p_{\text{up}}]$  by  $k_0$  samples. With  $[p_{\text{low}}, p_{\text{up}}]$  and  $\psi_{ISS}(\cdot)$ , we compute the input-specific sample size  $\hat{k}$ . Then we estimate the certified radius by sampling  $\hat{k}$  noisy samples. The algorithm is shown in Alg. 2.

**Compare ISS(AD) to IAS** We compare ISS to IAS in Fig. 3a, Fig. 3b where  $\hat{R}(k, p_A; \sigma) := \sigma \Phi^{-1}(\mathbf{B}(\alpha; p_A k, k - p_A k + 1))$ . As presented, IAS assigns 50,000 for both  $x_1, x_2$ , while ISS assigns 35,000 and 65,000 for  $x_1, x_2$  respectively. The sample size of ISS are computed by solving  $\text{AD}(k; 100,000, p_A) \leq 0.0042$ . For each certified radius, the decline in  $x_2$  certified radius is up to 0.0075 due to the sample size reduction  $100,000 \rightarrow 50,000$ , which is  $1.78 \times U_{AD}$  of ISS. For the average certified radius, ISS trades 0.002 radius of  $x_1$  for 0.003 radius of  $x_2$ , thus ACR of ISS improves IAS



(a)  $\hat{R}$ - $k$  curve of  $x_1$  ( $p_A = 0.51$ ). (b)  $\hat{R}$ - $k$  curve of  $x_2$  ( $p_A = 0.99$ ).

Figure 3: ISS certifies higher ACR on  $x_1, x_2$  than IAS.



(a)  $\hat{R}_{ISS} - \hat{R}_{IAS}$  w.r.t.  $p_A$  (b)  $p_A$  distribution (ImageNet).

Figure 4: ISS fits the practical smoothed classifiers.

0.0005 under the same average sample size. The improvement is because ISS tends to assign larger sample sizes to the high- $p_A$  inputs, which meets the property of  $\hat{R}(k, p_A; \sigma)$ . Namely,  $\hat{R}(k + \Delta k, p_A; \sigma) - \hat{R}(k, p_A; \sigma)$  increases with  $p_A$ , meaning that assigning larger sample sizes to the high- $p_A$  inputs is more efficient than input-agnostic sampling.

**ISS fits the well-trained smoothed classifiers** Fig. 4a reports  $\hat{R}_{ISS} - \hat{R}_{IAS}$  where  $\hat{R}_{ISS}$  denotes the radius certified by ISS of  $\bar{k} = 100,000$ ,  $U_{AD} = 0.05$  and  $\hat{R}_{IAS}$  denotes the radius certified by IAS at  $k = 30,000^2$ . We observe that ISS certifies higher certified radius when  $p_A > 0.94$ . Fig. 4b reports the  $p_A$  distribution of the test set<sup>3</sup> on the real ImageNet base classifier ( $\sigma = 0.5$ ) trained by Consistency (Jeong and Shin 2020). We found that the probability mass of  $p_A$  distribution is concentrated around  $p_A = 1.0$ , which is the interval where  $\hat{R}_{ISS} - \hat{R}_{IAS} > 0$ . Furthermore, ISS is expected to outperform IAS on the smoothed classifiers trained by other algorithms, since their  $p_A$  distributions have the similar property (see appendix).

## 5 Experiments

We evaluate our proposed method ISS on two benchmark datasets: CIFAR-10 (Krizhevsky 2009) and ImageNet (Rusakovsky et al. 2015). All the experiments are conducted on CPU (16 Intel(R) Xeon(R) Gold 5222 CPU @ 3.80GHz) and GPU (one NVIDIA RTX 2080 Ti). We observe that the certification runtime is roughly proportional to the average sample size when fixing the model architecture, as shown in

<sup>2</sup>Here we choose to compare ISS to IAS ( $k = 30,000$ ) is because that the average sample size of ISS ( $\bar{k} = 100,000$ ,  $U_{AD} = 0.05$ ) on the ImageNet model trained by Consistency (Jeong and Shin 2020) ( $\sigma = 0.5$ ) is roughly 30,000.

<sup>3</sup>We sample  $k = 1000,000$  Monte Carlo samples and approximately regard  $k_A/k$  as the exact value of  $p_A$ .

Table 1. The hyperparameters are listed in Table 2. For clarity,  $ISS - c_1 - c_2$  denotes ISS at  $\bar{k} = c_1 \cdot 10,000$ ,  $U_{AD} = c_2$ , and  $IAS - c_1$  denotes IAS at  $k = c_1 \cdot 10,000$ . The overhead of computing  $\psi_{ISS}$  is reported in Table 3.

### 5.1 Evaluation metrics

Our evaluation metrics include average sample size, runtime, MAD, ACR and certified accuracy, where MAD denotes the maximum absolute decline between the radius certified before and after the acceleration among all the testing data<sup>4</sup>. ACR and certified accuracy  $CA(r)$  at the radius  $r$  are computed as follows:

$$ACR := \frac{1}{|\mathcal{D}_{test}|} \sum_{(x,y) \in \mathcal{D}_{test}} R(x; g) \cdot \mathbb{I}(g(x) = y) \quad (7)$$

$$CA(r) := \frac{1}{|\mathcal{D}_{test}|} \sum_{(x,y) \in \mathcal{D}_{test}} \mathbb{I}(R(x; g) > r) \cdot \mathbb{I}(g(x) = y) \quad (8)$$

where  $R(x; g)$  denotes the estimated certified radius of  $g(x)$ .

### 5.2 Overall analysis of ACR and runtime

Fig. 5c, Fig. 5d, Fig. 5a, Fig. 5b present the overall empirical results of ISS and IAS on CIFAR-10 and ImageNet. As presented, ISS significantly accelerates the certification for randomized smoothing. Specifically, on ImageNet ( $\sigma = 0.5, 1.0$ ),  $ISS - 10 - 0.05$ ,  $ISS - 10 - 0.1$  reduce the original time cost 962 minutes (the green dotted lines) to roughly 300, 200 respectively at  $U_{AD} = 0.05, 0.1$  respectively. Overall, the speedups of ISS are even higher on CIFAR-10. We also compare ISS to IAS on two datasets. We found that ISS always achieves higher ACR than IAS in the similar time cost. For ImageNet ( $\sigma = 1.0$ ),  $ISS - 20 - 0.05$  even further improves IAS - 20 by a moderate margin, while the time cost of  $ISS - 20 - 0.05$  is only  $0.56 \times$  of IAS - 20. The full results are reported in the supplemental material.

### 5.3 Results of ISS (AD) on ImageNet

Table 1 reports the results of ISS<sup>5</sup>. Remarkably, ISS reduces the average sample size to roughly  $\frac{3}{10} \times$ ,  $\frac{1}{5} \times$  at the cost of  $U_{AD} = 0.05, 0.10$  respectively, meaning the speedups are roughly  $\frac{10}{3} \times$ ,  $5 \times$ . We found that the MADs of IAS are higher than ISS, meaning that IAS will cause a large radius decline on the specific inputs. Namely, the MAD of IAS - 10.4 is more than  $7 \times$   $ISS_{AD-50-0.05}$ . ISS consistently surpasses IAS on ACR.  $ISS_{AD-50-0.10}(\sigma = 1.0)$  achieves  $ACR = 1.005$  in 796 minutes while IAS only achieves  $ACR = 0.976$  in 1,000 minutes. We also observe that ISS slightly lower than IAS on the low-radius certified accuracies. It is because ISS tends to assign the small sample sizes to those inputs with low  $p_A$ , which inevitably sacrifices the certified radii of

<sup>4</sup>Note the speedup of ISS deterministically depends on the  $p_A$  distribution of the testing set. Since the smoothed classifiers trained by different training algorithms, including SmoothAdv (Salman et al. 2019), MACER (Zhai et al. 2020) and Consistency (Jeong and Shin 2020), report the similar  $p_A$  distributions, ISS will perform similarly on the models trained by other algorithms.

<sup>5</sup>Here we only report the results at  $\sigma = 0.5$ , and  $\sigma = 1.0$  because the work (Jeong and Shin 2020) only releases the training hyperparameters at  $\sigma = 0.5, 1.0$  for consistency training algorithm.



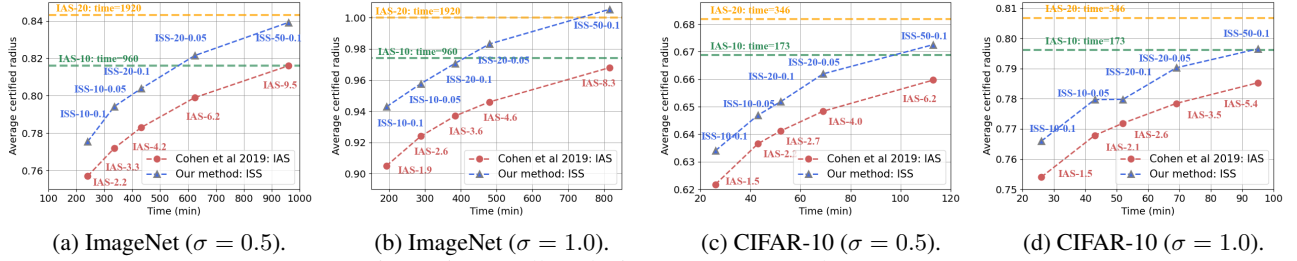


Figure 5: Overall analysis on ImageNet and CIFAR-10.

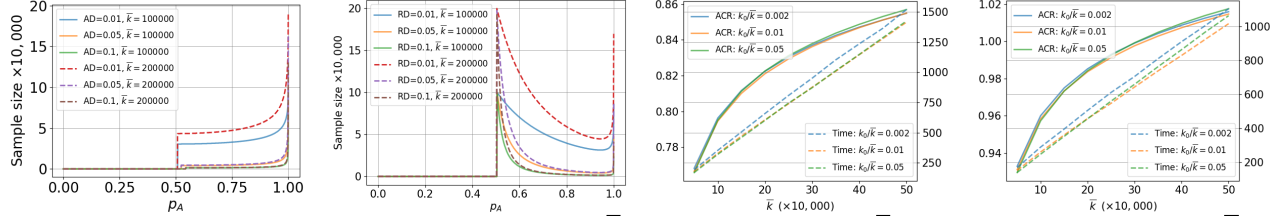


Figure 6: Ablation studies. **Upper:**  $k$ - $p_A$  curves w.r.t. AD,  $\bar{k}$  (Upper Left) and RD,  $\bar{k}$  (Upper Right). **Lower:** ACR- $\bar{k}$  curves and Time- $\bar{k}$  curves w.r.t.  $k_0/\bar{k}$  on ImageNet  $\sigma = 0.5$  (Lower Left) and ImageNet  $\sigma = 1.0$  (Lower Right).

Table 1: ImageNet: compare ISS to IAS on average sample size (Avg), certification runtime, maximum absolute decline (MAD), average certified radii (ACR) and certified accuracies (%) on the models trained by Consistency (Jeong and Shin 2020). Results are evaluated on 500 trials by  $id = 0, 100, \dots, 49800, 49900$ . The bold denotes better performance under the similar setting.

$\sigma$	Method	Avg	Time (min)	MAD	ACR	0.00	0.50	1.00	1.50	2.00	2.5	3.0	3.5	4.0
0.50	ISS <sub>AD-10-0.05</sub>	<b>32992</b>	<b>317</b>	<b>0.05</b>	<b>0.794</b>	54.6	49.8	42.4	<b>33.2</b>	0.0	0.0	0.0	0.0	0.0
	IAS - 3.3	33000	<b>317</b>	0.14	0.77	<b>54.8</b>	<b>50.2</b>	<b>43.4</b>	32.8	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-10-0.10</sub>	<b>22144</b>	<b>213</b>	<b>0.10</b>	<b>0.775</b>	54.6	49.8	42.0	<b>33.0</b>	0.0	0.0	0.0	0.0	0.0
	IAS - 2.2	22200	<b>213</b>	0.19	0.752	<b>54.8</b>	<b>50.2</b>	<b>43.4</b>	32.6	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-50-0.05</sub>	<b>144220</b>	<b>1385</b>	<b>0.05</b>	<b>0.856</b>	<b>54.8</b>	<b>50.2</b>	42.8	<b>34.2</b>	<b>29.8</b>	0.0	0.0	0.0	0.0
	IAS - 14.4	144400	1386	0.15	0.831	<b>54.8</b>	<b>50.2</b>	<b>43.4</b>	33.4	0.0	0.0	0.0	0.0	0.0
1.00	ISS <sub>AD-50-0.10</sub>	<b>95381</b>	<b>916</b>	<b>0.10</b>	<b>0.839</b>	<b>54.8</b>	<b>50.2</b>	42.8	<b>34.2</b>	0.0	0.0	0.0	0.0	0.0
	IAS - 9.5	95400	<b>916</b>	0.20	0.815	<b>54.8</b>	<b>50.2</b>	<b>43.4</b>	33.2	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-10-0.05</sub>	<b>25987</b>	<b>250</b>	<b>0.05</b>	<b>0.958</b>	40.6	36.8	31.8	<b>28.0</b>	<b>24.0</b>	<b>20.2</b>	<b>17.4</b>	<b>13.4</b>	0.0
	IAS - 2.6	26000	<b>250</b>	0.35	0.917	<b>41.0</b>	<b>37.2</b>	<b>32.6</b>	<b>28.0</b>	<b>24.0</b>	20.0	16.6	0.0	0.0
	ISS <sub>AD-10-0.10</sub>	<b>19209</b>	<b>185</b>	<b>0.10</b>	<b>0.943</b>	40.2	36.8	31.8	27.4	<b>24.0</b>	<b>20.0</b>	<b>17.4</b>	<b>13.4</b>	0.0
	IAS - 1.9	19400	186	0.43	0.903	<b>41.0</b>	<b>37.0</b>	<b>32.2</b>	<b>28.0</b>	<b>24.0</b>	<b>20.0</b>	16.2	0.0	0.0
	ISS <sub>AD-50-0.05</sub>	<b>104037</b>	<b>999</b>	<b>0.05</b>	<b>1.015</b>	40.6	36.8	32.2	<b>28.0</b>	<b>24.0</b>	<b>20.4</b>	<b>17.4</b>	<b>13.4</b>	<b>13.4</b>
	IAS - 10.4	104200	1000	0.37	0.976	<b>41.4</b>	<b>37.4</b>	<b>32.6</b>	<b>28.0</b>	<b>24.0</b>	<b>20.4</b>	<b>17.4</b>	<b>13.4</b>	0.0
	ISS <sub>AD-50-0.10</sub>	<b>82899</b>	<b>796</b>	<b>0.10</b>	<b>1.005</b>	40.6	36.8	32.2	<b>28.0</b>	<b>24.0</b>	20.0	<b>17.4</b>	<b>13.4</b>	<b>13.4</b>
	IAS - 8.3	83000	797	0.43	0.967	<b>41.4</b>	<b>37.4</b>	<b>32.6</b>	<b>28.0</b>	<b>24.0</b>	<b>20.4</b>	<b>17.4</b>	<b>13.4</b>	0.0

Table 2: Experiment setting.

Dataset	CIFAR-10	ImageNet
Model	ResNe-110	ResNet-50
Training by	MACER	Consistency
$\bar{k}$	100,000, 500,000	
$k_0$	0.01 $\bar{k}$	
$\sigma$	0.25, 0.5, 1.0	0.5, 1.0

Table 3: Runtime for computing  $\psi_{ISS}$ .

AD	Time (s)	RD	Time (s)
0.01	0.70	0.01	39.47
0.05	0.65	0.05	13.52
0.10	0.57	0.10	7.50

low- $p_A$  inputs. Meanwhile, ISS significantly improves the high-radius certified accuracies and ACR in return.

## 5.4 Results of ISS (RD) on ImageNet

Table 4 reports the results of ISS (RD) on ImageNet at  $U_{RD} = 0.05, 0.10$ . ISS reduces the average sample size to roughly  $\frac{7}{10} \times, \frac{7}{20} \times$  at controllable cost of RD = 1%, 5% respectively. Compared to IAS, ISS (RD) also improves ACR.

## 5.5 Results of ISS (AD) on CIFAR-10

Table 5 reports the results of ISS (OF AD) on CIFAR-10. ISS reduces the average sample size to roughly  $\frac{1}{5} \times, \frac{1}{10} \times$  at  $U_{AD} = 0.05, 0.10$ . Remarkably, ISS still improves ACRs and MADs, high-radius certified accuracies by a moderate margin on CIFAR-10. These empirical comparisons suggest that ISS is a better acceleration.

Table 4: ImageNet: comparison on average sample size (Avg), certification runtime (in minutes), average certified radii (ACR) and certified accuracies (%) on models trained by Consistency (Jeong and Shin 2020).

$\sigma$	Method	Avg	Time	ACR	0.00	0.50	1.00	1.50	2.00	2.5	3.0	3.5	4.0
0.50	ISS <sub>RD-10-0.01</sub>	<b>70919</b>	<b>682</b>	<b>0.809</b>	<b>54.8</b>	<b>50.2</b>	<b>43.4</b>	<b>33.2</b>	0.0	0.0	0.0	0.0	0.0
	IAS - 7.1	71000	<b>682</b>	0.803	<b>54.8</b>	<b>50.2</b>	<b>43.4</b>	<b>33.2</b>	0.0	0.0	0.0	0.0	0.0
	ISS <sub>RD-10-0.05</sub>	<b>34591</b>	333	<b>0.781</b>	<b>54.8</b>	<b>50.2</b>	42.2	<b>33.0</b>	0.0	0.0	0.0	0.0	0.0
	IAS - 3.5	34600	<b>332</b>	0.772	<b>54.8</b>	<b>50.2</b>	<b>43.4</b>	<b>32.8</b>	0.0	0.0	0.0	0.0	0.0
1.00	ISS <sub>RD-10-0.01</sub>	<b>67207</b>	<b>646</b>	<b>0.966</b>	<b>41.4</b>	<b>37.4</b>	<b>32.6</b>	<b>28.0</b>	<b>24.0</b>	<b>20.4</b>	<b>17.4</b>	<b>13.4</b>	0.0
	IAS - 6.7	67400	647	0.959	41.2	<b>37.4</b>	<b>32.6</b>	<b>28.0</b>	<b>24.0</b>	<b>20.4</b>	<b>17.4</b>	<b>13.4</b>	0.0
	ISS <sub>RD-10-0.05</sub>	<b>32515</b>	<b>313</b>	<b>0.935</b>	<b>41.2</b>	37.0	31.8	27.6	23.8	<b>20.0</b>	<b>17.2</b>	<b>13.4</b>	0.0
	IAS - 3.3	32600	<b>313</b>	0.927	41.0	<b>37.2</b>	<b>32.6</b>	<b>28.0</b>	<b>24.0</b>	<b>20.0</b>	16.8	<b>13.4</b>	0.0

Table 5: CIFAR-10: comparison on average sample size (Avg), certification runtime (in minutes), maximum absolute decline (MAD), average certified radii (ACR) and certified accuracies (%) on the models trained by MACER (Zhai et al. 2020). Results are evaluated on 500 testing data of  $id = 0, 20, \dots, 9960, 9980$ . The bold denotes better performance.

$\sigma$	Method	Avg	Time	MAD	ACR	0.00	0.25	0.50	0.75	1.00	1.25	1.50	1.75	2.00	2.25	2.50	2.75	3.0	3.25	3.5	3.75	4.00
0.25	ISS <sub>AD-10-0.05</sub>	<b>22237</b>	<b>39</b>	<b>0.05</b>	<b>0.492</b>	76.8	68.0	49.4	<b>38.8</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 2.2	22400	<b>39</b>	0.10	0.483	<b>77.8</b>	<b>68.6</b>	<b>52.0</b>	37.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-10-0.10</sub>	<b>10945</b>	<b>19</b>	<b>0.10</b>	<b>0.473</b>	76.8	68.0	48.8	<b>37.6</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 1.1	11000	<b>19</b>	0.15	0.462	<b>77.4</b>	<b>68.4</b>	<b>51.6</b>	36.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-50-0.05</sub>	<b>98984</b>	172	<b>0.05</b>	<b>0.529</b>	77.4	68.4	51.6	<b>39.8</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 9.9	99000	<b>171</b>	0.10	0.518	<b>77.8</b>	<b>69.0</b>	<b>52.2</b>	39.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-50-0.10</sub>	<b>46509</b>	<b>81</b>	<b>0.10</b>	<b>0.512</b>	77.4	68.4	51.6	<b>39.4</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 4.7	46600	<b>81</b>	0.14	0.501	<b>77.8</b>	<b>68.8</b>	<b>52.0</b>	38.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
0.50	ISS <sub>AD-10-0.05</sub>	<b>21836</b>	<b>38</b>	<b>0.05</b>	<b>0.647</b>	60.6	53.0	46.8	39.8	32.4	<b>26.0</b>	<b>19.8</b>	<b>13.0</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 2.2	22000	<b>38</b>	0.20	0.633	<b>61.8</b>	<b>54.0</b>	<b>47.8</b>	<b>40.2</b>	<b>32.8</b>	<b>26.0</b>	19.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-10-0.10</sub>	<b>14620</b>	<b>26</b>	<b>0.10</b>	<b>0.634</b>	60.6	53.0	46.8	39.4	31.0	<b>26.0</b>	<b>19.6</b>	<b>11.2</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 1.5	14800	<b>26</b>	0.25	0.621	<b>61.8</b>	<b>54.0</b>	<b>47.8</b>	<b>40.2</b>	<b>32.6</b>	<b>26.0</b>	19.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-50-0.05</sub>	<b>91293</b>	<b>158</b>	<b>0.05</b>	<b>0.68</b>	61.8	54.0	47.6	40.0	32.6	26.0	<b>20.2</b>	<b>14.2</b>	<b>10.2</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 9.1	91400	<b>158</b>	0.20	0.667	<b>62.2</b>	<b>54.4</b>	<b>48.0</b>	<b>40.2</b>	<b>33.0</b>	<b>26.6</b>	19.8	13.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	ISS <sub>AD-50-0.10</sub>	<b>61567</b>	<b>107</b>	<b>0.10</b>	<b>0.673</b>	61.8	54.0	47.6	40.0	32.6	<b>25.2</b>	<b>20.0</b>	<b>13.8</b>	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IAS - 6.2	61600	<b>107</b>	0.25	0.659	<b>62.2</b>	<b>54.4</b>	<b>47.8</b>	<b>40.2</b>	<b>33.0</b>	<b>26.4</b>	19.8	12.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1.00	ISS <sub>AD-10-0.05</sub>	<b>21153</b>	<b>37</b>	<b>0.05</b>	<b>0.78</b>	42.6	40.2	37.2	33.4	30.4	27.0	24.4	21.2	<b>18.4</b>	<b>14.6</b>	<b>13.4</b>	<b>10.4</b>	<b>8.8</b>	<b>6.2</b>	<b>4.0</b>	<b>3.2</b>	0.0
	IAS - 2.1	21200	<b>37</b>	0.40	0.763	<b>42.8</b>	<b>40.6</b>	<b>37.4</b>	<b>34.0</b>	<b>31.0</b>	<b>27.4</b>	<b>24.8</b>	<b>21.4</b>	<b>18.4</b>	<b>14.6</b>	12.8	9.8	8.0	4.6	0.0	0.0	0.0
	ISS <sub>AD-10-0.10</sub>	<b>14751</b>	<b>26</b>	<b>0.10</b>	<b>0.766</b>	42.4	39.6	36.8	32.8	30.0	26.4	24.0	20.6	18.2	<b>14.4</b>	<b>12.8</b>	<b>10.2</b>	<b>8.8</b>	<b>6.0</b>	<b>4.0</b>	0.0	0.0
	IAS - 1.5	14800	<b>26</b>	0.50	0.754	<b>42.8</b>	<b>40.4</b>	<b>37.4</b>	<b>33.8</b>	<b>30.8</b>	<b>27.4</b>	<b>24.8</b>	<b>21.4</b>	<b>18.4</b>	<b>14.4</b>	<b>12.8</b>	9.6	7.2	3.8	0.0	0.0	0.0
	ISS <sub>AD-50-0.05</sub>	<b>70204</b>	<b>123</b>	<b>0.05</b>	<b>0.803</b>	<b>42.8</b>	40.4	<b>37.4</b>	33.8	30.4	27.4	24.6	<b>21.4</b>	<b>18.4</b>	14.6	13.4	<b>10.4</b>	<b>9.2</b>	<b>6.8</b>	<b>4.8</b>	<b>4.0</b>	<b>3.2</b>
	IAS - 7.0	70400	<b>123</b>	0.47	0.79	<b>42.8</b>	<b>40.6</b>	<b>37.4</b>	<b>34.4</b>	<b>31.0</b>	<b>27.6</b>	<b>25.0</b>	<b>21.4</b>	<b>18.4</b>	<b>14.8</b>	<b>13.6</b>	10.2	8.8	6.0	4.0	0.0	0.0
	ISS <sub>AD-50-0.10</sub>	<b>54102</b>	<b>94</b>	<b>0.10</b>	<b>0.797</b>	<b>42.8</b>	40.4	<b>37.4</b>	33.8	30.4	27.4	24.6	21.2	<b>18.4</b>	14.4	13.0	<b>10.0</b>	<b>9.2</b>	<b>6.6</b>	<b>4.8</b>	<b>4.0</b>	<b>3.2</b>
	IAS - 5.4	54200	<b>94</b>	0.53	0.785	<b>42.8</b>	<b>40.6</b>	<b>37.4</b>	<b>34.4</b>	<b>31.0</b>	<b>27.6</b>	<b>25.0</b>	<b>21.4</b>	<b>18.4</b>	<b>14.8</b>	<b>13.4</b>	<b>10.0</b>	8.6	5.6	4.0	0.0	0.0

## 5.6 Ablation study

**Choice on AD or RD** As shown in Fig. 6, when  $p_A : p_A \in [0.5, 1.0]$  increases, the sample size of ISS (AD) monotonically increases, while the sample size of ISS (RD) first decreases and then increases around  $p_A = 1.0$ . ISS (AD) can greatly improve ACR, but tends to sacrifice the certified radii of low- $p_A$  inputs a relatively larger proportion. ISS (RD) sacrifices all inputs the same proportion of radius.

**Impact of  $p_A$  and  $\bar{k}$**  We investigate the impact of  $p_A$  and  $\bar{k}$  in Fig. 6 (Upper). For both AD and RD, the sample size is 0 when  $p_A \leq 0.5$ . It is because the certified radius is 0 when  $p_A \leq 0.5$ . As expected, the sample size monotonically increases with  $\bar{k}$  and decreases with AD (or RD).

**Impact of  $k_0/\bar{k}$**  We investigate the impact  $k_0/\bar{k}$  on the runtime and ACR in Fig. 6 (Lower). Too small  $k_0/\bar{k}$  results in a loose confidence interval  $[p_{\text{low}}, p_{\text{up}}]$ , which can cause the ISS sample size  $\hat{k}$  to be much larger than required. Too large  $k_0/\bar{k}$  may waste too much computation in estimating  $[p_{\text{low}}, p_{\text{up}}]$ . Our choice  $k_0/\bar{k} = 0.01$  performs well across various noise levels on CIFAR-10 and ImageNet.

## 6 Conclusion

Randomized smoothing has been suffering from the long certification runtime, but the current acceleration methods are low-efficiency. Therefore, we propose input-specific sampling, which adaptively assigns the sample size. Our work establishes an initial step towards a better performance-time trade-off for the certification of randomized smoothing. Specifically, Our strong empirical results suggest that ISS is a promising acceleration. Specifically, ISS speeds up the certification by more than  $4\times$  only at the controllable cost of 0.10 certified radius on ImageNet. An interesting direction for future work is to make the confidence interval estimation method adapt to the input.

## Acknowledgements

This work was supported by the National Key Research and Development Program of China No. 2020YFB1806700, Shanghai Municipal Science and Technology Major Project Grant 2021SHZDZX0102, NSFC Grant 61932014, NSFC Grant 61872241, Project BE2020026, the Key R&D Program of Jiangsu, China. This work is also partially sup-

ported by The Hong Kong Polytechnic University under Grant P0030419, P0030929, and P0035358.

## References

- Athalye, A.; Carlini, N.; and Wagner, D. 2018. Obfuscated gradients give a false sense of security: circumventing defenses to adversarial examples. In *ICML*.
- Awasthi, P.; Jain, H.; Rawat, A. S.; and Vijayaraghavan, A. 2020. Adversarial robustness via robust low rank representations. In *NeurIPS*.
- Buckman, J.; Roy, A.; Raffel, C.; and Goodfellow, I. 2018. Thermometer encoding: One hot way to resist adversarial examples. In *ICLR*.
- Cisse, M.; Bojanowski, P.; Grave, E.; Dauphin, Y.; and Usunier, N. 2017. Parseval networks: Improving robustness to adversarial examples. In *ICML*.
- Clopper, C. J.; and Pearson, E. S. 1934. The use of confidence or fiducial limits illustrated in the case of the binomial. *Biometrika*.
- Cohen, J. M.; Rosenfeld, E.; and Kolter, J. Z. 2019. Certified Adversarial Robustness via Randomized Smoothing. In Chaudhuri, K.; and Salakhutdinov, R., eds., *ICML*.
- Dhillon, G. S.; Azzadenesheli, K.; Lipton, Z. C.; Bernstein, J.; Kossai, J.; Khanna, A.; and Anandkumar, A. 2018. Stochastic activation pruning for robust adversarial defense. In *ICLR*.
- Eykholt, K.; Evtimov, I.; Fernandes, E.; Li, B.; Rahmati, A.; Tramer, F.; Prakash, A.; Kohno, T.; and Song, D. 2018. Physical adversarial examples for object detectors. *USENIX Workshop*.
- Feng, H.; Wu, C.; Chen, G.; Zhang, W.; and Ning, Y. 2020. Regularized training and tight certification for randomized smoothed classifier with provable robustness. In *AAAI*.
- Fischer, M.; Baader, M.; and Vechev, M. 2020. Certified Defense to Image Transformations via Randomized Smoothing. In *NeurIPS*.
- Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2015. Explaining and harnessing adversarial examples. In *ICLR*.
- Guo, C.; Rana, M.; Cisse, M.; and van der Maaten, L. 2018. Countering Adversarial Images using Input Transformations. In *ICLR*.
- Jeong, J.; and Shin, J. 2020. Consistency regularization for certified robustness of smoothed classifiers. In *NeurIPS*.
- Jia, J.; Cao, X.; Wang, B.; and Gong, N. Z. 2020. Certified robustness for top-k predictions against adversarial perturbations via randomized smoothing. In *ICLR*.
- Jia, J.; and Gong, N. 2018. AttriGuard: A practical defense against attribute inference attacks via adversarial machine learning. In *USENIX Security Symposium*.
- Krizhevsky, A. 2009. Learning multiple layers of features from tiny images. Technical report, Department of Computer Science, University of Toronto.
- Kurakin, A.; Goodfellow, I. J.; and Bengio, S. 2017. Adversarial machine learning at scale. In *ICLR*.
- Levine, A.; and Feizi, S. 2020. Robustness Certificates for sparse adversarial attacks by randomized ablation. In *AAAI*.
- Li, L.; Weber, M.; Xu, X.; Rimanic, L.; Xie, T.; Zhang, C.; and Li, B. 2020. Provable robust learning based on transformation-specific smoothing. *arXiv*.
- Ma, X.; Li, B.; Wang, Y.; Erfani, S. M.; Wijewickrema, S.; Schoenebeck, G.; Song, D.; Houle, M. E.; and Bailey, J. 2018. Characterizing adversarial subspaces using local intrinsic dimensionality. In *ICLR*.
- Raghunathan, A.; Steinhardt, J.; and Liang, P. S. 2018. Semidefinite relaxations for certifying robustness to adversarial examples. In *NeurIPS*.
- Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A. C.; and Fei-Fei, L. 2015. ImageNet large scale visual recognition challenge. *IJCV*.
- Salman, H.; Li, J.; Razenshteyn, I.; Zhang, P.; Zhang, H.; Bubeck, S.; and Yang, G. 2019. Provably robust deep learning via adversarially trained smoothed classifiers. In *NeurIPS*.
- Song, Y.; Kim, T.; Nowozin, S.; Ermon, S.; and Kushman, N. 2018. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *ICLR*.
- Svoboda, J.; Masci, J.; Monti, F.; Bronstein, M. M.; and Guibas, L. 2019. Peernets: Exploiting peer wisdom against adversarial attacks. In *ICLR*.
- Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; and Fergus, R. 2014. Intriguing properties of neural networks. In *ICLR*.
- Wong, E.; and Kolter, J. Z. 2018. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *ICML*.
- Wong, E.; Schmidt, F.; Metzen, J. H.; and Kolter, J. Z. 2018. Scaling provable adversarial defenses. In *NeurIPS*.
- Xie, C.; Wang, J.; Zhang, Z.; Ren, Z.; and Yuille, A. 2017. Mitigating adversarial effects through randomization. *arXiv*.
- Yang, G.; Duan, T.; Hu, E.; Salman, H.; Razenshteyn, I. P.; and Li, J. 2020. Randomized smoothing of all shapes and sizes. In *ICML*.
- Zhai, R.; Dan, C.; He, D.; Zhang, H.; Gong, B.; Ravikumar, P.; Hsieh, C.; and Wang, L. 2020. MACER: attack-free and scalable robust training via maximizing certified Radius. In *ICLR*.
- Zhang, H.; Chen, H.; Xiao, C.; Li, B.; Boning, D.; and Hsieh, C.-J. 2019. Towards Stable and Efficient Training of Verifiably Robust Neural Networks. *arXiv*.