

# CERTIFIED $\ell_2$ ATTRIBUTION ROBUSTNESS VIA UNIFORMLY SMOOTHED ATTRIBUTIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Model attribution is a popular tool to explain the rationales behind model predictions. However, recent work suggests that the attributions are vulnerable to minute perturbations, which can be added to input samples to manipulate the attributions while maintaining the prediction outputs. Although empirical studies have shown positive performance via adversarial training, an effective certified defense method is eminently needed to understand the robustness of attributions. In this work, we propose to use uniform smoothing technique that augments the vanilla attributions by noises uniformly sampled from a certain space. It is proved that, for all perturbations within the attack region, the cosine similarity between uniformly smoothed attribution of perturbed sample and the unperturbed sample is guaranteed to be lower bounded. We also derive alternative formulations of the certification that is equivalent to the original one and provides the maximum size of perturbation or the minimum smoothing radius such that the attribution can not be perturbed. We evaluate the proposed method on three datasets and show that the proposed method can effectively protect the attributions from attacks, regardless of the architecture of networks, training schemes and the size of the datasets.

## 1 INTRODUCTION

The developments and wider uses of deep learning models in various security-sensitive applications, such as autonomous driving, medical diagnosis, and legal judgments, have raised discussions of the trustworthiness of these models. The lack of explainability of deep learning models, especially the recent popular large language models (Brown et al., 2020), has been one of the main concerns. Regulators have started to require the explainability of the AI models in some applications (Goodman & Flaxman, 2017), and the explainability of AI models has been one of the main focuses of the research community (Doshi-Velez & Kim, 2017). Model attributions, as one of the important tools to explain the rationales behind the model predictions, have been used to understand the decision-making process of the models. For example, medical practitioners use the explanations generated by attribution methods to assist them in making important medical decisions (Antoniadi et al., 2021; Hrinivich et al., 2023; Du et al., 2022). Similarly, in autonomous vehicles, the explainability helps to deal with potential liability and responsibility gaps (Atakishiyev et al., 2021; Burton et al., 2020) and people naturally requires the confirmation of the safety critical decisions. However, since these users often lack expertise in machine learning and technical details of attributions, there is a risk that the attributions have been manipulated without noticing. Attackers may specifically target attributions to mislead investigators, propagate false narratives, or evade detection, potentially leading to serious consequences. Consequently, practitioners could lose trust in these methods, resulting in their refusal to use these methods. Thus, a trustworthy application requires not only the predictions made by the AI models, but also its explainability produced from attributions. However, the attributions have also been shown recently to be vulnerable to small perturbations (Ghorbani et al., 2019). Similar to adversarial attacks, attribution attacks generate perturbations that can be added to input samples. These perturbations distort the attributions while maintaining unchanged prediction outputs. This misleadingly gives a false sense of security to practitioners who blindly trust the attributions. An effective defense method is emergently needed to protect the attributions from the attacks.

Unlike adversarial defense, which has been extensively investigated to mitigate the harm of adversarial attacks that using both empirical (Madry et al., 2018; Athalye et al., 2018; Carlini & Wagner,

2017) and certified (Cohen et al., 2019; Wong & Kolter, 2018; Yang et al., 2020; Lecuyer et al., 2019) defense methods, attribution defense is neglected. Almost all attribution protection works focus on adversarially training the model by augmenting the training data with manipulated samples to improve the robustness of attributions (Boopathy et al., 2020; Chen et al., 2019; Ivankay et al., 2020; Sarkar et al., 2021; Wang & Kong, 2022). Levine et al. (2019) and a series of subsequent works (Liu et al., 2022; Huai et al., 2022; Gu et al., 2023) study the certifications of attributions under categorical ranking measurements, which are not easy to extend to other domains. A recent work by Wang & Kong (2023) attempts to derive a practical upper bound of cosine similarity for the worst-case attribution deviation, while suffering from strict assumptions and heavy computations. Overall, there is a gap in scalable methods for providing generalized certification of attribution robustness. In this work, we put our focus on the smoothed version of attributions and seek to provide a theoretical guarantee that the attributions are robust to any type of perturbations within  $\ell_2$  attack budget.

Based on the previously defined formulation of attribution robustness (Wang & Kong, 2023), given a network  $f$ , its attribution function  $g$  and perturbation  $\delta \in \mathbb{R}^d$ , the attribution robustness is defined as the optimal value that maximizes the worst-case attribution difference  $D(\cdot, \cdot)$  under provided attack budget,

$$\max_{\delta} D(g(\mathbf{x}), g(\mathbf{x} + \delta)) \quad \text{s.t.} \quad \|\delta\| \leq \epsilon. \quad (1)$$

However, the aforementioned study only provides an approximate method to solve the optimization problem under the strict assumptions that the networks is locally linear, and, as a result, is unable to generalize to modern neural networks. To give a complete certification on the problem, we follow the formulation and attempt to find the effective upper bound of the attribution difference, equivalently, the lower bound of attribution similarity, which can be applied to any network. We propose to use uniformly smoothed attribution, which is a smoothed version of the original attribution, and show that, for all perturbations within the allowable attack budget, the cosine similarity that measures the difference between perturbed and unperturbed uniformly smoothed attribution can be certified to be lower bounded. The contribution of this paper can be summarized as follows:

- We provide a theoretical guarantee that demonstrates the robustness of the uniformly smoothed attribution to any perturbations within allowable region. The robustness is measured by the similarity between the perturbed and unperturbed smoothed attribution. The method can be generally applied to any neural networks, and can be efficiently scaled to larger size images. To the best knowledge of the authors, this is the first work that provides a theoretical guarantee for attribution robustness.
- We present alternative formulations of the certification that are equivalent to the original one and also practical to be implemented. The alternative formulations determine the maximum size of perturbation, or the minimum radius of smoothing, ensuring that the attribution remains within a given tolerance.
- We demonstrate that the uniform smoothing can protect the attribution against  $\ell_2$  attacks. More importantly, we evaluated the proposed method on the well-bounded integrated gradients and show that it can be effectively implemented and can successfully protect and certify the attributions from  $\ell_2$  attacks.

The rest of this paper is organized as follows. In Section 2, we review the related works. In Section 3 and Section 4, we introduce the uniformly smoothed attribution and show that it can be certified against attribution attacks. In Section 5, we present experimental results and evaluate the proposed method. Finally, we conclude this paper and in Section 6.

## 2 RELATED WORKS

### 2.1 ATTRIBUTION METHODS

Attribution methods study the importance of each input feature, and measure that how much every feature contributes to the model prediction. One of the most popular attribution approaches is the gradient-based method. Based on the property that gradient is the measurement of the rate of change, the gradient-based methods measure the feature importance by weighting the gradients in different

ways. Examples of gradient-based methods include saliency (Simonyan et al., 2014), integrated gradients (IG) (Sundararajan et al., 2017), full-gradient (Srinivas & Fleuret, 2019), and *etc.* Other attribution methods include occlusion (Simonyan et al., 2014), which measures the importance of each input feature by occluding the feature and measuring the change of the model prediction, layer-wise relevance propagation (LRP) (Bach et al., 2015), which propagates the output relevance to the input layer, and SHAP related methods (Lundberg & Lee, 2017; Sundararajan & Najmi, 2020; Kwon & Zou, 2022). An important property of many attribution methods is the axiom of completeness that  $\sum_i g_i(\mathbf{x}) = f_j(\mathbf{x})$ , which indicates the relationship between attribution and the prediction score. Note that the gradient-based attribution methods are upper-bounded since the gradients of the model output with respect to the input are upper-bounded.

## 2.2 ATTRIBUTION ATTACKS AND DEFENSES

Ghorbani et al. (2019) first pointed out that attributions can be fragile to iterative attribution attack, and Dombrowski et al. (2019) extended the attack to be targeted that attributions can be changed purposely into any preset patterns. Similar to adversarial attacks, attribution attacks maximize the loss function that measures the difference between the original attributions and the target attributions. In addition, the attribution attacks are controlled not to alter the classification results. To defend against attribution attacks, adversarial training (Madry et al., 2018) approaches have been adopted. Chen et al. (2019) and Boopathy et al. (2020) minimize the  $\ell_p$ -norm differences between perturbed and original attributions, and Ivankay et al. (2020) considers Pearson’s correlation coefficient. It is worth noting that these methods empirically improve attribution robustness. Meanwhile, a series of certification methods Levine et al. (2019); Liu et al. (2022); Huai et al. (2022); Gu et al. (2023) have been proposed to ensure that attribution changes do not exceed a certain threshold under any perturbations within the allowable attack region. These methods measure attribution changes using *top-k intersection*, while a continuous alternative, cosine similarity, remains unexplored. It has been proved that using cosine similarity to measure attribution differences is consistent as top-k intersection and Kendall’s rank correlation (Wang & Kong, 2022), and it is more likely to extend to other domains.

## 2.3 RANDOMIZED SMOOTHING

The smoothing technique has been popular in improving certified adversarial robustness (Liu et al., 2018; Lecuyer et al., 2019; Cohen et al., 2019; Yang et al., 2020). The smoothed classifiers take a batch of inputs that are randomly sampled from the neighbourhood of original inputs under certain distributions  $\mu$  and make the decisions based on the most likely outputs, *i.e.*,  $\arg \max_y \mathbb{P}_{\eta \sim \mu}[F(\mathbf{x} + \boldsymbol{\eta}) = y]$ . They provide the certification of the a radius such that no perturbation within the radius can alter the classification label. Cohen et al. (2019) certifies the  $\ell_2$  attack based on the Neyman-Pearson lemma and the result is alternatively proved by Salman et al. (2019) using explicit Lipschitz constants. Lecuyer et al. (2019) and Teng et al. (2020) consider Laplacian smoothing for the  $\ell_1$  attack. Yang et al. (2020) derived a similar result in  $\ell_\infty$  though the radius becomes small when the dimension of data gets large. Kumar & Goldstein (2021) applied randomized smoothing to structured output, which the attributions belong to, but the specific bound for attribution is too loose to be meaningful. Thus, there are no existing works that provide defense and valid certifications for attribution using randomized smoothing due to the difficulty of defining the attribution robustness and the computation of the attribution gradient. In this work, an effective method to formulate the smoothed attribution robustness as a simple optimization problem is proposed to provide certifications against attribution attacks.

## 3 UNIFORMLY SMOOTHED ATTRIBUTION

Consider a classifier  $f : \mathbb{R}^d \rightarrow [0, 1]^c$  that maps the input  $\mathbf{x} \in \mathbb{R}^d$  to the softmax output  $y \in [0, 1]$ , and its attribution function  $g(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . The *smoothed attribution* of  $f$  is to construct a new attribution  $h$  by taking the mean of attributions on  $\mathbf{x} + \boldsymbol{\eta}$ , where  $\boldsymbol{\eta}$  is randomly drawn from a density  $\mu$ , *i.e.*, the smoothed attribution  $h$  can be defined as follows:

$$h(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\eta} \sim \mu}[g(\mathbf{x} + \boldsymbol{\eta})]. \tag{2}$$

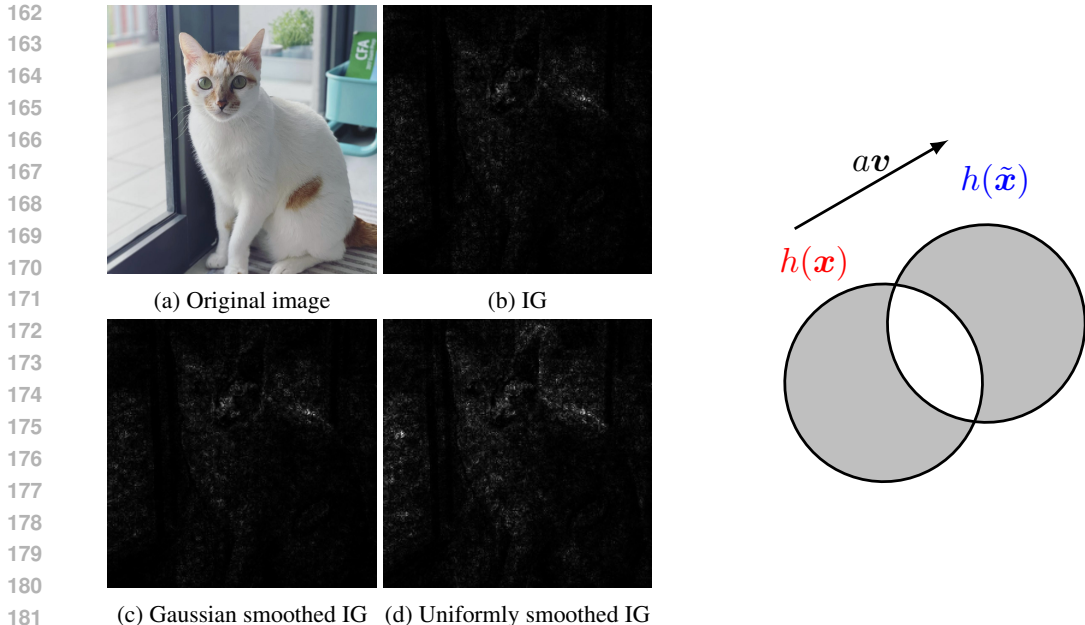


Figure 1: **(left)** Examples of attributions (zoom in for better visibility). We choose to show the integrated gradients (IG) and its corresponding smoothing results. For Gaussian smoothing, the noise level is set to  $\sigma = 0.2$  and for uniformly smoothed IG,  $\ell_2$  ball with radius  $\sqrt{3}\sigma$  is used. **(right)** A 2D illustration of the volumes of  $\mathcal{B}(\mathbf{x}; r)$  and  $\mathcal{B}(\tilde{\mathbf{x}}; r)$ , as well as the relationship between  $h(\mathbf{x})$  and  $h(\tilde{\mathbf{x}})$ . Here  $h(\mathbf{x})$  is the original attribution, and  $av$  represents the magnitude and direction of the translation of  $h(\mathbf{x})$  after the sample is perturbed.  $V_U$  in Theorem 1 is the volume of shaded region in the figure, and  $V_S$  is the volume of each individual ball. When  $h(\mathbf{x})$  is fixed, the lower bound of the cosine similarity between  $h(\mathbf{x})$  and  $h(\mathbf{x}) + av$  can be derived as a function of volumes

To construct smoothed attribution, the density  $\mu$  can be chosen arbitrarily, and the smoothed attributions provide visually sharpened gradient-based attributions (see Figure 1 (left)). Smilkov et al. (2017) choose  $\mu$  to be multivariate Gaussian distribution on input gradients to weaken the visually noisy attribution. In this work, we aim at certifying the attribution robustness under  $\ell_2$  attack; thus we choose  $\mu$  to be a uniform distribution on a  $d$ -dimensional closed space  $\mathcal{S}$  centered at  $\mathbf{0}$ , especially the  $\ell_2$ -norm ball of radius  $r$ ,  $\mathcal{B}(\mathbf{0}; r) = \{\mathbf{y} : \|\mathbf{y}\|_2 \leq r, \mathbf{y} \in \mathbb{R}^d\}$ . It can be seen that smoothing under this setting also provides high attribution quality (Figure 1d). To quantitatively evaluate the effectiveness of uniformly smoothed attribution, we can further evaluate its performance using GridPG introduced by Rao et al. (2022), which quantifies the significance of individual features in terms of positive contributions or influences. The GridPG values of IG, uniformly smoothed IG and Gaussian smoothed IG for 5,000 randomly selected ImageNet examples are 0.4021, 0.4093 and 0.4110, respectively, which suggest that the uniformly smoothed attributions can achieve comparable performance of GridPG with the Gaussian smoothed attributions, as well as the original non-smoothed attributions.

#### 4 CERTIFYING THE COSINE SIMILARITY OF SMOOTHED ATTRIBUTIONS

We now consider  $\mathcal{S}$  as an  $\ell_2$ -norm ball for the ease of analyzing the certification. Adapted from the formulation in Eq. (1), the robustness of attribution is defined as the minimum possible attribution similarity when a natural image is perturbed by attribution attacks. As mentioned in Section 2, this work studies cosine similarity as the measurement of similarity, as it has been shown to be the most suitable alternative to the non-differentiable Kendall’s rank correlation, the most common evaluation index (Wang & Kong, 2022). Suppose that the  $\ell_2$  attribution attack is performed upon input sample  $\mathbf{x}$ , the maximum allowable perturbation is  $\epsilon$ , i.e.,  $\tilde{\mathbf{x}} = \mathbf{x} + \delta$ , where  $\|\delta\|_2 \leq \epsilon$ . For all  $\|\delta\|_2 \leq \epsilon$ , we want to find out the minimum value of  $\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}}))$ , i.e., the lower bound of cosine similarity

for attributions. Thus, the optimization problem we are interested in is formulated as follows:

$$\min_{\delta} \cos(h(\mathbf{x}), h(\mathbf{x} + \delta)) \quad \text{s.t.} \quad \|\delta\|_2 \leq \epsilon. \quad (3)$$

That is, we will show that, given an arbitrary sample point  $\mathbf{x}$ , for all perturbed samples  $\tilde{\mathbf{x}}$ , the cosine similarity between the original smoothed attribution and the corresponding perturbed smoothed attributions is guaranteed to be lower bounded by the optimum of (3). Alternatively, in a more practical perspective, given a threshold  $T$  for the cosine similarity, we want to know the maximum size of perturbation, or the minimum smoothing radius, such that no perturbations inside would cause the attribution difference to exceed the threshold.

However, although optimizing the cosine similarity for attribution can be an intuitive way to find the lower bound, it is difficult in this problem to directly study the cosine function. Moreover, it is also intractable to optimize the cosine similarity with respect to a vector  $\delta$ . To address this issue, we first reformulate the problem into an optimization over two scalars and then solve the alternative problem to obtain the lower bound of cosine similarity. All the proofs and derivations of theorems, lemmas and corollaries are provided in the Appendix A.

#### 4.1 ONE-DIMENSIONAL REFORMULATION

We note that cosine similarity of attributions is in fact an inner product of their normalized vectors. Besides, we also observe that  $h(\mathbf{x})$  is the mean of  $g(\mathbf{x} + \boldsymbol{\eta})$  with respect to  $\boldsymbol{\eta}$ , which is, by definition, equivalent to the integral of  $g$  weighted by the density of uniform distribution over  $\mathcal{B}(\mathbf{x}; r)$ . More importantly, for perturbed example  $\tilde{\mathbf{x}}$ , the weighting region is  $\mathcal{B}(\tilde{\mathbf{x}}; r)$ , which is a translated version of  $\mathcal{B}(\mathbf{x}; r)$  and they are expected to intersect with each other when the distance of their centers is smaller than twice of the radius  $r$ . Therefore, given the input sample  $\mathbf{x}$  and its attribution  $h(\mathbf{x})$ , we can rewrite the minimization of cosine similarity as follows:

$$\min_{a, \mathbf{v}} \frac{h(\mathbf{x})^T}{\|h(\mathbf{x})\|} \left( \frac{h(\mathbf{x}) + a\mathbf{v}}{\|h(\mathbf{x}) + a\mathbf{v}\|} \right) \quad (4)$$

where  $a$  represents the magnitude of the translation and the unit vector  $\mathbf{v}$  is the direction of translation as shown in Figure 1 (right). It can be shown that the magnitude  $a$  is constrained by a constant related to the intersecting volume and the property of the attribution function itself.

In a high-dimensional case, for a fixed cosine similarity value, a given  $h(\mathbf{x})$  and  $h(\tilde{\mathbf{x}})$ , in fact, form a spherical cone. Thus, we can decompose the directional unit vector  $\mathbf{v}$  into  $\mathbf{v} = \cos\theta\mathbf{v}_{\parallel} + \sin\theta\mathbf{v}_{\perp}$ , where  $\mathbf{v}_{\perp}$  is perpendicular to  $h(\mathbf{x})$  and  $\mathbf{v}_{\parallel}$  is parallel to  $h(\mathbf{x})$ . Then, we have

$$\min_{a, \theta} \frac{\|h(\mathbf{x})\| + a \cos\theta}{\sqrt{(\|h(\mathbf{x})\| + a \cos\theta)^2 + (a \sin\theta)^2}}, \quad (5)$$

where  $0 \leq \theta \leq 2\pi$ .

Since  $a$  is a scalar representing the magnitude of the translation from  $h(\mathbf{x})$  to  $h(\tilde{\mathbf{x}})$ , it can be shown that the magnitude of  $a$  is upper bounded by the magnitude of the gradient weighted by the ratio of volume change during the translation process. Specifically,  $a \leq MV_U/V_S$ , where  $V_S$  is the volume of the  $\ell_2$ -ball  $\mathcal{B}(\mathbf{0}; r)$ ,  $V_U$  is the volume of the union of the two sampling space centered at  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$  minus their intersection, and  $M$  is a constant that depends on the upper bound of  $g$ . We notice that the gradient-based attribution is a function of input gradient,  $\nabla f(\mathbf{x})$ , which is bounded for Lipschitz continuous networks (See Lemma 1 in the Appendix A); thus, the upper bound can be derived separately for different attribution functions.

#### 4.2 LOWER BOUND OF COSINE SIMILARITY FOR BOUNDED ATTRIBUTIONS

Now that we have reformulated the optimization with respect to vector  $\mathbf{v}$  into an alternative simpler one with respect to scalar values  $a$  and  $\theta$ . By solving the alternative problem, our result shows that the smoothed attribution is robust within the following half-angle of a spherical cone.

**Theorem 1.** *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be an upper bounded attribution function, and  $\boldsymbol{\eta} \stackrel{U}{\sim} \mathcal{B}(\mathbf{0}; r)$ . Let  $h$  be the smoothed version of  $g$  as defined in (2). Then, for all  $\tilde{\mathbf{x}} \in \{\mathbf{x} + \delta \mid \|\delta\|_2 \leq \epsilon\}$ , we have*

Table 1: Comparison of top-k and Kendall’s rank correlation between  $\ell_2$  perturbed and non-perturbed attributions on standard and robust models using smoothed and non-smoothed attributions.

		Standard	IG-NORM (Chen et al., 2019)	TRADES (Zhang et al., 2019)	IGR (Wang & Kong, 2022)
SM	top-k	0.3449	0.6575	0.6450	0.8354
	Kendall	0.1496	0.4709	0.4642	0.7553
SmoothSM	top-k	0.3853	0.6261	0.6082	0.8363
	Kendall	0.1670	0.4238	0.4119	0.7568
IG	top-k	0.4742	0.7075	0.6821	0.8402
	Kendall	0.1744	0.5098	0.5030	0.7839
SmoothIG	top-k	0.5302	0.6730	0.6528	0.8460
	Kendall	0.3819	0.4494	0.4533	0.7612

$\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}})) \geq T$ , where

$$T = \frac{\|h(\mathbf{x})\|_2}{\sqrt{\|h(\mathbf{x})\|_2^2 + M^2 V_U^2 / V_S^2}} \quad (6)$$

Here,  $M$  is the upper bound of  $g$ .  $V_S$  is the volume of the  $\ell_2$ -ball  $\mathcal{B}(\mathbf{0}; r)$ , and  $V_U$  is the volume of the union of the two sampling space centered at  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$  minus their intersection.

The entire proof can be found in Appendix A. The theorem points out that the lower bound of cosine similarity is related to the smoothing space around the input samples and the maximum allowable perturbation size. Moreover, when the smoothing space is a  $\ell_2$ -norm ball, the above result can be derived by directly computing two volumes  $V_U$  and  $V_S$ , which can be explicitly calculated by

$$V_U = 2V_S \times \left( 1 - I_{(2rh-h^2)/r^2} \left( \frac{d+1}{2}, \frac{1}{2} \right) \right) \quad (7)$$

where  $h = r - \epsilon/2 \geq 0$  and  $I_x(a, b)$  is the regularized incomplete beta function, cumulative density function of beta distribution (Li, 2010).

We observe the following properties of the above theorem.

1. Unlike previous attribution robustness works, such as Dombrowski et al. (2019), Singh et al. (2020), Boopathy et al. (2020) and Wang & Kong (2022), which require the networks to be twice-differentiable and need to change the ReLU activation into Softplus, the proposed result does not assume anything on the classifiers. Thus, it can be safely applied to any neural network and any architectures.
2. The lower bound depends on the radius of smoothing,  $r$ . The bound becomes larger as the radius of smoothing grows. In extreme cases when  $r$  tends to infinity, the smoothing spaces of two samples completely overlap, which corresponds to the same smoothed attribution and their cosine similarity becomes 1.
3. At the same time, the lower bound also decreases when the attack budget  $\epsilon$  increases. Besides, when  $\epsilon$  increases beyond the constraint that  $h = r - \epsilon/2 \geq 0$  and tends to infinity, the distance between two attributions will become further and their cosine similarity will tend to 0 in high dimensional space, which makes the lower bound becomes trivial.
4. The proposed method can be scaled to datasets with large images. The lower bound is efficient to compute since only the smoothed attribution of given sample needed. On the contrary, the previous works that approximately estimate the attribution robustness (Wang & Kong, 2023) require the computation of input Hessian and the corresponding eigenvalues and eigenvectors, which becomes intractable for larger size images on modern neural networks.

Table 2: The theoretical lower bound ( $T$  in Eqn. (6)) for cosine similarity evaluated on baseline models using MNIST. Note that the bound is not achievable for  $r = 0.5$  when  $\epsilon = 1.0$ , since the radius must be greater than  $\epsilon/2$ .

$\epsilon = 0.5$	$\ell_2$ radius ( $r$ )	0.5	1.0	1.5	2.0	2.5	3.0	3.5
	Standard	0.3002	0.3141	0.3385	0.3732	0.4144	0.4600	0.5057
	IG-NORM	0.4038	0.4189	0.4432	0.4729	0.5055	0.5466	0.5909
	IGR	0.4145	0.4269	0.4482	0.4792	0.5208	0.5748	0.6392
$\epsilon = 1.0$	$\ell_2$ radius ( $r$ )	0.5	1.0	1.5	2.0	2.5	3.0	3.5
	Standard	/	0.3034	0.3092	0.3264	0.3716	0.3990	0.4178
	IG-NORM	/	0.3650	0.3822	0.3892	0.4220	0.4974	0.5365
	IGR	/	0.3834	0.4025	0.4558	0.4914	0.5237	0.5358

### 4.3 ALTERNATIVE FORMULATIONS OF THE ATTRIBUTION ROBUSTNESS

The previous section formulates the robustness of smoothed attribution in terms of the smallest cosine similarity between the original and perturbed smoothed attribution, when the attack budget and the smoothing radius are fixed. In some scenarios, practitioners want to formulate the robustness in different ways. For example, we may want to find the maximum allowable perturbation  $\epsilon$  such that the cosine similarity between the original and perturbed smoothed attribution is guaranteed to be greater than a predefined threshold  $T$ . On the other hand, one can also obtain the minimum smoothing radius needed such that the desired attribution robustness is achieved within allowable attack region. The following corollary provides the alternative formulations of the attribution robustness.

**Corollary 1.** *Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a bounded attribution function, and  $\eta \stackrel{U}{\sim} \mathcal{B}(\mathbf{x}; r)$ . Let  $h$  be the smoothed version of  $g$  as defined in (2).*

- (i) *Given a predefined threshold  $T \in [0, 1]$ , then for all  $\|\delta\|_2 \leq \epsilon$ , we have  $\cos(h(\mathbf{x}), h(\mathbf{x} + \delta)) \geq T$ , where*

$$\epsilon = 2r \sqrt{1 - I_Z^{-1} \left( \frac{d+1}{2}, \frac{1}{2} \right)}. \quad (8)$$

- (ii) *Given a predefined threshold  $T \in [0, 1]$  and the maximum perturbation size  $\epsilon \geq 0$ , the smoothed attribution satisfies  $\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}})) \geq T$  for all  $\tilde{\mathbf{x}} \in \{\mathbf{x} + \delta \mid \|\delta\|_2 \leq \epsilon\}$  when  $r \geq R$ , where*

$$R = \frac{\epsilon}{2} \left( 1 - I_Z^{-1} \left( \frac{d+1}{2}, \frac{1}{2} \right) \right)^{-\frac{1}{2}}. \quad (9)$$

$I_Z^{-1}(a, b)$  is the inverse of the regularized incomplete beta function, and  $Z$  is defined as

$$Z = 1 - \frac{\|h(\mathbf{x})\|_2}{2M} \left( \frac{1}{T^2} - 1 \right) \quad (10)$$

The derivation of the Corollary can be found in Appendix A. We notice that, in these formulations, when the smoothing radius is larger, the maximum allowable perturbation is also larger, which allows stronger attacks while keeping the attribution similarity within a controllable range. Similarly, when the maximum allowable perturbation is larger, the minimum smoothing radius is also larger, which means that the attribution similarity can be maintained with a larger smoothing radius.

## 5 EXPERIMENTS AND RESULTS

In this section, we evaluate the effectiveness of uniformly smoothed attribution. Following previous work on attribution robustness, we use the  $\ell_2$  attribution attack adapted from the Iterative Feature Importance Attacks (IFIA) by Ghorbani et al. (2019). It is first shown that the uniformly smoothed

Table 3: Theoretical lower bounds evaluated on non-robust model and ImageNet using different radii  $r$  and attack sizes  $\epsilon$ . Note that the method is only applicable when radius  $r$  must be greater than  $\epsilon/2$ .

$r$	0.5	1.0	1.5	2.0	2.5	3.0	3.5
$\epsilon = 0.5$	0.2612	0.2594	0.2704	0.2716	0.2892	0.2973	0.3029
$\epsilon = 1.0$	/	0.1803	0.1992	0.1746	0.1904	0.2127	0.2502
$\epsilon = 2.0$	/	/	0.1753	0.1852	0.2044	0.2015	0.2045

attributions are less likely to be perturbed comparing with non-smoothed attributions when being attacked by IFIA. After which, we present the results of certification using the proposed lower bound of cosine similarity. The experiments are conducted on baseline models including adversarial robust models (Madry et al., 2018; Zhang et al., 2019), attributional robust models (Chen et al., 2019; Singh et al., 2020; Ivankay et al., 2020; Wang & Kong, 2022), as well as non-robust models trained for standard classification tasks. Following those baseline models, the method is tested on the validation sets of MNIST (LeCun et al., 2010) using a small-size convolutional network, on CIFAR-10 (Krizhevsky, 2009) using a ResNet-18 (He et al., 2016), and on ImageNet (Russakovsky et al., 2015) using a ResNet-50 (He et al., 2016). More details of the experiments are described in the Appendix. All experiments are run on NVIDIA GeForce RTX 3090.<sup>1</sup>

To empirically compute the uniformly smoothed attribution for every sample,  $N$  points are randomly sampled from the  $d$ -dimensional sphere uniformly, and augmented to the input sample. To do so, the sampling technique introduced by Box & Muller (1958) is applied. The integration in Eqn. 2 is then empirically estimated using Monte Carlo integration where the computation scales linearly with the number of samples, *i.e.*,  $\hat{h}(\mathbf{x}) = \frac{1}{N} \sum g(\mathbf{x} + \boldsymbol{\eta}_i)$ , for  $\boldsymbol{\eta}_i \stackrel{U}{\sim} \mathcal{B}(\mathbf{0}; r)$ . For large  $N$ , the estimator  $\hat{h}(\mathbf{x})$  almost surely converges to  $h(\mathbf{x})$  (Feller, 1991); hence the convergence of  $\hat{T}$  to  $T$  can be obtained (see Appendix C.1 for details). Unless specifically stated, we choose the number of samples  $N$  to be 100,000 to compute the proposed lower bound, and  $N^* = 300$  for the uniformly smoothed attribution being attacked in all experiments.

### 5.1 EVALUATION OF THE ROBUSTNESS OF UNIFORMLY SMOOTHED ATTRIBUTION

We first conduct the experiment to verify that the uniformly smoothed attribution itself is more robust than the original attribution. The uniform smoothing around the  $\ell_2$  ball with radius 0.5 is applied to the saliency map (SM) (Simonyan et al., 2014) and integrated gradients (IG) (Sundararajan et al., 2017) and evaluate on CIFAR-10. The resultant attributions are denoted by SmoothSM and SmoothIG, respectively. We then attack the attributions using the  $\ell_2$  IFIA attack and evaluate the robustness using Kendall’s rank correlation and top-k intersection (Ghorbani et al., 2019). The experiments are evaluated on both non-robust model (*Standard*) and robust models (*IG-NORM*, *TRADES*, *IGR*). Note that IFIA is directly performed on the smoothSM and smoothIG, instead of its original counterpart. Since the PGD-like attribution attack requires to take the derivative of the attribution to determine the direction of gradient descent, the double backpropagation is needed. Thus, it is necessary to replace the ReLU activation by the twice-differentiable Softplus during attack (Dombrowski et al., 2019). The results are shown in Table 1.

We observe that for the non-robust model, both SmoothSM and SmoothIG perform better than its non-smoothed counterparts in both metrics, which shows that the uniformly smoothed attribution itself is more resistant to the attribution attacks. For models that are specifically trained to defend against the attribution attacks using heuristic methods adapted from adversarial training, *e.g.*, IG-NORM, TRADES and IGR, the smoothed attributions show comparable robustness to the non-smoothed attribution. Moreover, we also notice that SmoothIG performs better than SmoothSM, especially for the non-robust models. This can be attributed to the fact that IG satisfies the axiom of completeness, which ensures that the sum of IG is upper-bounded by the model output. In addition, we also observed that the smoothing technique does not always enhance the robustness of attribution, as measured by top-k and Kendall’s rank correlation. It is worth noting that Yeh et al. (2019) argued that randomized smoothing can reduce attribution sensitivities and consequently improve ro-

<sup>1</sup>Source code will be released later.



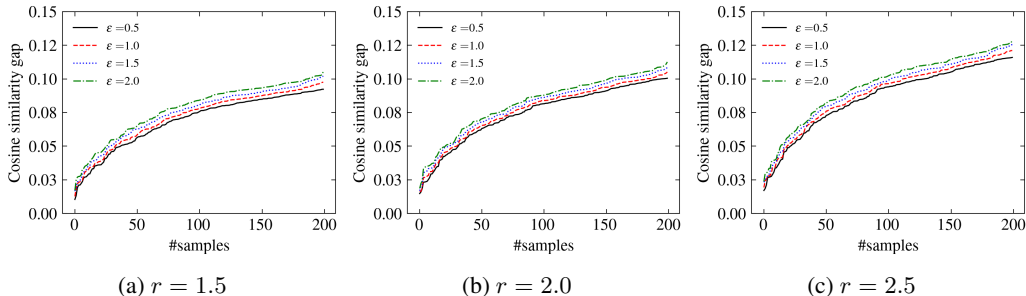


Figure 2: The gap between theoretical bounds and empirical cosine similarity between original and perturbed attribution evaluated on CIFAR-10 using IGR.

bustness. The disparity in our findings stems from the fact that we specifically evaluated Kendall’s rank correlation, whereas their study defined sensitivity based on  $\ell_2$ -norm distance, which was subsequently deemed inappropriate for evaluating attribution robustness by Wang & Kong (2022).

## 5.2 EVALUATION OF THE CERTIFICATION OF UNIFORMLY SMOOTHED ATTRIBUTIONS

In this section, the lower bound of the cosine similarity between the original and attacked attribution is reported. We use the integrated gradients as an example since it is well-bounded due to the axiom of completeness, and the technique can be also applied to any other gradient-based attributions.

Table 2 reports the theoretical bound evaluated on MNIST computed using Theorem 1. We include the non-robust and two attributional robust models and compute the bound for different  $\ell_2$  radius  $r$  and attack size  $\epsilon$  pairs. The lower bounds are validated by examining the actual  $\ell_2$  attacks. Specifically, each input sample has been attacked 20 times and the cosine similarities of resulting perturbed attributions with original attributions are examined. In total, 200,000 attacked images are tested for each parameter pair and none of the evaluation metrics exceeds the theoretical bound. Moreover, we also observe that the bound becomes tighter when the  $\ell_2$  radius  $r$  increases and when the attack size  $\epsilon$  decreases. Besides, since IGR is more robust than IG-NORM (Wang & Kong, 2022), we can observe that the lower bound is also a valid measurement of the robustness of the models.

In Figure 2, we show the gap between the theoretical lower bound and the empirical cosine similarity between the original and perturbed attributions. The results are evaluated on CIFAR-10 using IGR. Out of 10,000 testing samples, the 200 with the smallest gaps are chosen for each pair of  $r$  and  $\epsilon$ , and the gaps are sorted for better visualization. We notice that the gaps between the theoretical bound and empirical cosine similarity are positive and small, which shows the validity and the tightness of the proposed bound.

In Table 3, we also include the theoretical lower bound evaluated on ImageNet to show that the proposed method is also applicable to large-scale datasets. Since the current attribution attacks and attribution defense methods do not scale to large-scale datasets, we only include the non-robust model. Since our method does not rely on the second-order derivative of the output with respect to the input, it can be scaled to ImageNet-size datasets. For the experiments on ResNet-50, each certification for one single sample takes around 15 seconds. We observe that the reported bounds are also consistent with our theoretical findings.

## 6 CONCLUSION

In this paper, we attempt to use the uniformly smoothed attribution to certify the attribution robustness evaluated by cosine similarity. The smoothed attribution is constructed by taking the mean of the attributions computed from input samples augmented by noises uniformly sampled from an  $\ell_2$  ball. It is proved that the cosine similarity between the original and perturbed smoothed attribution is lower-bounded based on a geometric formulation related to the volume of the hyperspherical cap. Alternative formulations are provided to find the maximum allowable size of perturbations and the minimum radius of smoothing in order to maintain the attribution robustness. The method works

on bounded gradient-based attribution methods for all convolutional neural networks and is scalable to large datasets. We empirically demonstrate that the method can be used to certify the attribution robustness, using the well-bounded integrated gradients, and the state-of-the-art attributional robust models on MNIST, CIFAR-10 and ImageNet.

## 7 LIMITATIONS AND BROADER IMPACTS

The method in paper can be generally applied to any convolutional neural networks and any bounded attribution methods. Although the existence of an upper bound has been shown for all gradient-based methods, in some extreme cases when the upper bound for certain attribution is trivial, *i.e.*, an extremely large value, the proposed lower bound for attribution robustness also becomes trivial. In future work, we will investigate the upper bound for other bounded attribution methods and provide the corresponding lower bounds for attribution robustness. Besides, our current smoothing technique is restricted to the uniform distribution, and we will explore other distributions for the smoothing technique in future work.

Our work attempts to draw the attention of the community to the need for a guarantee of attribution robustness. With the increasingly large number of applications of deep learning, the transparency and trustworthiness of neural networks are crucial for users to understand the outcomes and to avoid any abuse of the techniques. While the study of the security of networks could reveal their potential risks that can be misused, we believe this work has more positive impacts to the community and can encourage the development of more trustworthy deep learning applications.

## REFERENCES

- Anna Markella Antoniadis, Yuhan Du, Yasmine Guendouz, Lan Wei, Claudia Mazo, Brett A Becker, and Catherine Mooney. Current challenges and future opportunities for xai in machine learning-based clinical decision support systems: a systematic review. *Applied Sciences*, 11(11):5088, 2021.
- Shahin Atakishiyev, Mohammad Salameh, Hengshuai Yao, and Randy Goebel. Explainable artificial intelligence for autonomous driving: A comprehensive overview and field guide for future research directions. *arXiv preprint arXiv:2112.11561*, 2021.
- Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pp. 274–283. PMLR, 2018.
- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.
- Akhilan Boopathy, Sijia Liu, Gaoyuan Zhang, Cynthia Liu, Pin-Yu Chen, Shiyu Chang, and Luca Daniel. Proper network interpretability helps adversarial robustness in classification. In *International Conference on Machine Learning*, pp. 1014–1023. PMLR, 2020.
- George EP Box and Mervin E Muller. A note on the generation of random normal deviates. *The annals of mathematical statistics*, 29(2):610–611, 1958.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Simon Burton, Ibrahim Habli, Tom Lawton, John McDermid, Phillip Morgan, and Zoe Porter. Mind the gaps: Assuring the safety of autonomous systems from an engineering, ethical, and legal perspective. *Artificial Intelligence*, 279:103201, 2020.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57. IEEE Computer Society, 2017.
- Jiefeng Chen, Xi Wu, Vaibhav Rastogi, Yingyu Liang, and Somesh Jha. Robust attribution regularization. In *Advances in Neural Information Processing Systems*, pp. 14300–14310, 2019.

- 540 Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized  
541 smoothing. In *International Conference on Machine Learning*, pp. 1310–1320. PMLR, 2019.
- 542  
543 Abhranil Das and Wilson S Geisler. A method to integrate and classify normal distributions. *Journal*  
544 *of Vision*, 21(10):1–1, 2021.
- 545 Robert B Davies. The distribution of a linear combination of  $\chi^2$  random variables. *Journal of the*  
546 *Royal Statistical Society Series C: Applied Statistics*, 29(3):323–333, 1980.
- 547  
548 Ann-Kathrin Dombrowski, Maximillian Alber, Christopher Anders, Marcel Ackermann, Klaus-  
549 Robert Müller, and Pan Kessel. Explanations can be manipulated and geometry is to blame.  
550 In *Advances in Neural Information Processing Systems*, pp. 13589–13600, 2019.
- 551 Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning.  
552 *arXiv preprint arXiv:1702.08608*, 2017.
- 553  
554 Yuhan Du, Anthony R Rafferty, Fionnuala M McAuliffe, Lan Wei, and Catherine Mooney. An  
555 explainable machine learning-based clinical decision support system for prediction of gestational  
556 diabetes mellitus. *Scientific Reports*, 12(1):1170, 2022.
- 557 Pierre Duchesne and Pierre Lafaye De Micheaux. Computing the distribution of quadratic forms:  
558 Further comparisons between the liu–tang–zhang approximation and exact methods. *Computational*  
559 *Statistics & Data Analysis*, 54(4):858–862, 2010.
- 560 William Feller. *An introduction to probability theory and its applications, Volume 2*, volume 81.  
561 John Wiley & Sons, 1991.
- 562  
563 Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In  
564 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3681–3688, 2019.
- 565  
566 Bryce Goodman and Seth Flaxman. European union regulations on algorithmic decision-making  
567 and a “right to explanation”. *AI magazine*, 38(3):50–57, 2017.
- 568 Alex Gu, Tsui-Wei Weng, Pin-Yu Chen, Sijia Liu, and Luca Daniel. Certified interpretability ro-  
569 bustness for class activation mapping. *arXiv preprint arXiv:2301.11324*, 2023.
- 570  
571 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
572 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
573 770–778, 2016.
- 574 William Thomas Hrinivich, Tonghe Wang, and Chunhao Wang. Interpretable and explainable ma-  
575 chine learning models in oncology. *Frontiers in Oncology*, 13:1184428, 2023.
- 576  
577 Mengdi Huai, Jinduo Liu, Chenglin Miao, Liuyi Yao, and Aidong Zhang. Towards automating  
578 model explanations with certified robustness guarantees. In *Proceedings of the AAAI Conference*  
579 *on Artificial Intelligence*, volume 36, pp. 6935–6943, 2022.
- 580 Adam Ivankay, Ivan Girardi, Chiara Marchiori, and Pascal Frossard. Far: A general framework for  
581 attributional robustness. *arXiv preprint arXiv:2010.07393*, 2020.
- 582  
583 Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 584  
585 Aounon Kumar and Tom Goldstein. Center smoothing: Certified robustness for networks with  
586 structured outputs. *Advances in Neural Information Processing Systems*, 34:5560–5575, 2021.
- 587  
588 Yongchan Kwon and James Y Zou. Weightedshap: analyzing and improving shapley based feature  
589 attributions. *Advances in Neural Information Processing Systems*, 35:34363–34376, 2022.
- 590  
591 Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*.  
592 Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 593  
594 Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified  
595 robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security*  
596 *and Privacy (SP)*, pp. 656–672. IEEE, 2019.

- 594 Alexander Levine, Sahil Singla, and Soheil Feizi. Certifiably robust interpretation in deep learning.  
595 *arXiv preprint arXiv:1905.12105*, 2019.
- 596 Shengqiao Li. Concise formulas for the area and volume of a hyperspherical cap. *Asian Journal of*  
597 *Mathematics & Statistics*, 4(1):66–70, 2010.
- 599 Ao Liu, Xiaoyu Chen, Sijia Liu, Lirong Xia, and Chuang Gan. Certifiably robust interpretation via  
600 rényi differential privacy. *Artificial Intelligence*, 313:103787, 2022.
- 601 Xuanqing Liu, Minhao Cheng, Huan Zhang, and Cho-Jui Hsieh. Towards robust neural networks via  
602 random self-ensemble. In *Proceedings of the European Conference on Computer Vision (ECCV)*,  
603 pp. 369–385, 2018.
- 605 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances*  
606 *in neural information processing systems*, 30, 2017.
- 607 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. To-  
608 wards deep learning models resistant to adversarial attacks. In *International Conference on Learn-*  
609 *ing Representations*, 2018. URL <https://openreview.net/forum?id=rJzIBfZAb>.
- 611 Remigijus Paulavičius and Julius Žilinskas. Analysis of different norms and corresponding lipschitz  
612 constants for global optimization. *Technological and Economic Development of Economy*, 12(4):  
613 301–306, 2006.
- 614 Sukrut Rao, Moritz Böhle, and Bernt Schiele. Towards better understanding attribution methods.  
615 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
616 10223–10232, 2022.
- 617 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
618 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual  
619 recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 621 Hadi Salman, Jerry Li, Ilya Razenshteyn, Pengchuan Zhang, Huan Zhang, Sebastien Bubeck, and  
622 Greg Yang. Provably robust deep learning via adversarially trained smoothed classifiers. *Ad-*  
623 *vances in Neural Information Processing Systems*, 32, 2019.
- 624 Anindya Sarkar, Anirban Sarkar, and Vineeth N Balasubramanian. Enhanced regularizers for attri-  
625 butional robustness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35,  
626 pp. 2532–2540, 2021.
- 627 Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:  
628 Visualising image classification models and saliency maps. In *2nd International Conference on*  
629 *Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Workshop Track*  
630 *Proceedings*, 2014.
- 631 Mayank Singh, Nupur Kumari, Puneet Mangla, Abhishek Sinha, Vineeth N Balasubramanian, and  
632 Balaji Krishnamurthy. Attributional robustness training using input-gradient spatial alignment. In  
633 *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020,*  
634 *Proceedings, Part XXVII 16*, pp. 515–533. Springer, 2020.
- 636 Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad:  
637 removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- 638 Suraj Srinivas and François Fleuret. Full-gradient representation for neural network visualization.  
639 *Advances in neural information processing systems*, 32, 2019.
- 640 Mukund Sundararajan and Amir Najmi. The many shapley values for model explanation. In *Inter-*  
641 *national conference on machine learning*, pp. 9269–9278. PMLR, 2020.
- 643 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
644 *International Conference on Machine Learning*, pp. 3319–3328. PMLR, 2017.
- 645 Jiaye Teng, Guang-He Lee, and Yang Yuan.  $\ell_1$  adversarial robustness certificates: a  
646 randomized smoothing approach, 2020. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=H1lQIgrFDS)  
647 [H1lQIgrFDS](https://openreview.net/forum?id=H1lQIgrFDS).

648 Fan Wang and Adams Wai-Kin Kong. Exploiting the relationship between kendall’s rank correla-  
649 tion and cosine similarity for attribution protection. *Advances in Neural Information Processing*  
650 *Systems*, 35:20580–20591, 2022.

651 Fan Wang and Adams Wai-Kin Kong. A practical upper bound for the worst-case attribution devia-  
652 tions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
653 pp. 24616–24625, 2023.

654 Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer  
655 adversarial polytope. In *International Conference on Machine Learning*, pp. 5286–5295. PMLR,  
656 2018.

657 Greg Yang, Tony Duan, J Edward Hu, Hadi Salman, Ilya Razenshteyn, and Jerry Li. Randomized  
658 smoothing of all shapes and sizes. In *International Conference on Machine Learning*, pp. 10693–  
659 10705. PMLR, 2020.

660 Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the  
661 (in) fidelity and sensitivity of explanations. *Advances in Neural Information Processing Systems*,  
662 32, 2019.

663 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.  
664 Theoretically principled trade-off between robustness and accuracy. In *International Conference*  
665 *on Machine Learning*, pp. 7472–7482. PMLR, 2019.

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

## A PROOFS

**Lemma 1.** (Paulavičius & Žilinskas (2006)) For  $L$ -Lipschitz function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_q \quad (11)$$

where  $L = \max \{\|\nabla f(\mathbf{x})\|_p : \mathbf{x} \in S\}$  is Lipschitz constant. Thus,  $\|\nabla f(\mathbf{x})\|_p \leq L$ .

*Proof.* Refer to Paulavičius & Žilinskas (2006) for the proof.  $\square$

**Theorem 1.** Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a upper bounded attribution function, and  $\boldsymbol{\eta} \stackrel{U}{\sim} \mathcal{B}(\mathbf{0}; r)$ . Let  $h$  be the smoothed version of  $g$  as defined in (2). Then, for all  $\tilde{\mathbf{x}} \in \{\mathbf{x} + \boldsymbol{\delta} \mid \|\boldsymbol{\delta}\|_2 \leq \epsilon\}$ , we have  $\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}})) \geq T$ , where

$$T = \frac{\|h(\mathbf{x})\|_2}{\sqrt{\|h(\mathbf{x})\|_2^2 + M^2 V_U^2 / V_S^2}} \quad (6)$$

Here,  $M$  is the upper bound of  $g$ .  $V_S$  is the volume of the  $\ell_2$ -ball  $\mathcal{B}(\mathbf{0}; r)$ , and  $V_U$  is the volume of the union of the two sampling space centered at  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$  minus their intersection.

*Proof.* As defined in Eqn. (2)

$$h(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\eta} \sim \mathcal{B}(\mathbf{0}; r)} [g(\mathbf{x} + \boldsymbol{\eta})] = \frac{1}{V_S} \int_{\boldsymbol{\eta} \sim \mathcal{B}(\mathbf{0}; r)} g(\mathbf{x} + \boldsymbol{\eta}) d\boldsymbol{\eta} \quad (12)$$

where  $V_S$  is the volume of the  $\ell_p$ -ball with radius  $r$ . Similarly, let  $\tilde{\mathbf{x}} = \mathbf{x} + \boldsymbol{\delta}$ , where  $\boldsymbol{\delta} \in \mathbb{R}^d$  is a vector and  $\|\boldsymbol{\delta}\|_2 \leq \epsilon$ . Then, we have

$$h(\tilde{\mathbf{x}}) = \frac{1}{V_S} \int_{\boldsymbol{\eta} \sim \mathcal{B}(\mathbf{0}; r)} g(\tilde{\mathbf{x}} + \boldsymbol{\eta}) d\boldsymbol{\eta} \quad (13)$$

We note that when  $\boldsymbol{\eta} \sim \mathcal{B}(\mathbf{0}; r)$ ,  $\mathbf{x} + \boldsymbol{\eta} \sim \mathcal{B}(\mathbf{x}; r)$  and  $\tilde{\mathbf{x}} + \boldsymbol{\eta} \sim \mathcal{B}(\tilde{\mathbf{x}}; r)$ . We then rewrite  $h(\mathbf{x})$  and  $h(\tilde{\mathbf{x}})$  as follows:

$$h(\mathbf{x}) = \underbrace{\frac{1}{V_S} \int_{\mathbf{x} \sim \mathcal{B}(\mathbf{x}; r) \setminus \mathcal{B}(\tilde{\mathbf{x}}; r)} g(\mathbf{x}) d\mathbf{x}}_{R_1} + \underbrace{\frac{1}{V_S} \int_{\mathbf{x} \sim \mathcal{B}(\tilde{\mathbf{x}}; r) \cap \mathcal{B}(\mathbf{x}; r)} g(\mathbf{x}) d\mathbf{x}}_{R_2} \quad (14)$$

and

$$h(\tilde{\mathbf{x}}) = \underbrace{\frac{1}{V_S} \int_{\mathbf{x} \sim \mathcal{B}(\tilde{\mathbf{x}}; r) \cap \mathcal{B}(\mathbf{x}; r)} g(\mathbf{x}) d\mathbf{x}}_{R_2} + \underbrace{\frac{1}{V_S} \int_{\mathbf{x} \sim \mathcal{B}(\tilde{\mathbf{x}}; r) \setminus \mathcal{B}(\mathbf{x}; r)} g(\mathbf{x}) d\mathbf{x}}_{R_3} \quad (15)$$

Hence,

$$h(\tilde{\mathbf{x}}) = h(\mathbf{x}) - R_1 + R_3 \quad (16)$$

Denote  $av = R_3 - R_1$ , where  $v$  is a unit vector in the same direction of  $R_3 - R_1$  and  $a = \|R_3 - R_1\|_2$  is a scalar with the same magnitude of  $R_3 - R_1$ . Then, we have

$$\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}})) = \frac{h(\mathbf{x})^\top}{\|h(\mathbf{x})\|_2} \left( \frac{h(\mathbf{x}) + av}{\|h(\mathbf{x}) + av\|_2} \right) \quad (17)$$

Note that the attribution  $g(\mathbf{x})$  is upper bounded by  $M$ , specifically,  $\|g(\mathbf{x})\|_2 \leq M$ , for some constant  $M$ . Thus, we can derive that

$$a = \|R_3 - R_1\|_2 \quad (18)$$

$$= \left\| \frac{1}{V_S} \left( \int_{\mathbf{x} \sim \mathcal{B}(\tilde{\mathbf{x}}; r) \setminus \mathcal{B}(\mathbf{x}; r)} g(\mathbf{x}) d\mathbf{x} - \int_{\mathbf{x} \sim \mathcal{B}(\mathbf{x}; r) \setminus \mathcal{B}(\tilde{\mathbf{x}}; r)} g(\mathbf{x}) d\mathbf{x} \right) \right\|_2 \quad (19)$$

$$\leq \frac{1}{V_S} \left( \left\| \int_{\mathbf{x} \sim \mathcal{B}(\tilde{\mathbf{x}}; r) \setminus \mathcal{B}(\mathbf{x}; r)} g(\mathbf{x}) d\mathbf{x} \right\|_2 + \left\| \int_{\mathbf{x} \sim \mathcal{B}(\mathbf{x}; r) \setminus \mathcal{B}(\tilde{\mathbf{x}}; r)} g(\mathbf{x}) d\mathbf{x} \right\|_2 \right) \quad (20)$$

$$\leq \frac{1}{V_S} \left( \int_{\mathbf{x} \sim \mathcal{B}(\tilde{\mathbf{x}}; r) \setminus \mathcal{B}(\mathbf{x}; r)} \|g(\mathbf{x})\|_2 d\mathbf{x} + \int_{\mathbf{x} \sim \mathcal{B}(\mathbf{x}; r) \setminus \mathcal{B}(\tilde{\mathbf{x}}; r)} \|g(\mathbf{x})\|_2 d\mathbf{x} \right) \quad (21)$$

$$\leq \frac{1}{V_S} \left( \int_{\mathbf{x} \sim \mathcal{B}(\tilde{\mathbf{x}}; r) \setminus \mathcal{B}(\mathbf{x}; r)} M d\mathbf{x} + \int_{\mathbf{x} \sim \mathcal{B}(\mathbf{x}; r) \setminus \mathcal{B}(\tilde{\mathbf{x}}; r)} M d\mathbf{x} \right) \quad (22)$$

$$= M \times \frac{V_{\mathcal{B}(\mathbf{x}; r) \setminus \mathcal{B}(\tilde{\mathbf{x}}; r)} \cup V_{\mathcal{B}(\tilde{\mathbf{x}}; r) \setminus \mathcal{B}(\mathbf{x}; r)}}{V_S} = M \frac{V_U}{V_S} \quad (23)$$

Thus, the lower bound of  $\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}}))$  can be found by solving the optimization problem <sup>2</sup>

$$\begin{aligned} \min_{\mathbf{v}} \quad & \frac{h(\mathbf{x})^\top}{\|h(\mathbf{x})\|} \left( \frac{h(\mathbf{x}) + a\mathbf{v}}{\|h(\mathbf{x}) + a\mathbf{v}\|} \right) \\ \text{s.t.} \quad & \|\mathbf{v}\| = 1 \\ & a \leq M \frac{V_U}{V_S} \end{aligned} \quad (24)$$

Since  $h(\mathbf{x})$  and  $h(\tilde{\mathbf{x}})$  form a spherical cone, we can decompose  $\mathbf{v}$  by  $\mathbf{v} = \cos \theta \mathbf{v}_\parallel + \sin \theta \mathbf{v}_\perp$ , where  $\mathbf{v}_\parallel$  and  $\mathbf{v}_\perp$  are two orthogonal unit vectors such that  $h^\top(\mathbf{x})\mathbf{v}_\perp = 0$  and  $\mathbf{v}_\parallel = h(\mathbf{x})/\|h(\mathbf{x})\|$ . Then, the optimization problem can be rewritten as

$$\min \quad \mathbf{v}_\parallel^\top \left( \frac{h(\mathbf{x}) + a(\cos \theta \mathbf{v}_\parallel + \sin \theta \mathbf{v}_\perp)}{\|h(\mathbf{x}) + a(\cos \theta \mathbf{v}_\parallel + \sin \theta \mathbf{v}_\perp)\|} \right) \quad (25)$$

$$\Rightarrow \min \quad \mathbf{v}_\parallel^\top \left( \frac{\|h(\mathbf{x})\| \mathbf{v}_\parallel + a(\cos \theta \mathbf{v}_\parallel + \sin \theta \mathbf{v}_\perp)}{\| \|h(\mathbf{x})\| \mathbf{v}_\parallel + a(\cos \theta \mathbf{v}_\parallel + \sin \theta \mathbf{v}_\perp) \|} \right) \quad (26)$$

$$\Rightarrow \min \quad \frac{(\|h(\mathbf{x})\| + a \cos \theta) \mathbf{v}_\parallel^\top \mathbf{v}_\parallel + a \sin \theta \mathbf{v}_\parallel^\top \mathbf{v}_\perp}{\sqrt{(\|h(\mathbf{x})\| + a \cos \theta)^2 \mathbf{v}_\parallel^\top \mathbf{v}_\parallel + (a \sin \theta)^2 \mathbf{v}_\perp^\top \mathbf{v}_\perp}} \quad (27)$$

$$\Rightarrow \min \quad \frac{\|h(\mathbf{x})\| + a \cos \theta}{\sqrt{(\|h(\mathbf{x})\| + a \cos \theta)^2 + (a \sin \theta)^2}} \quad (28)$$

Since  $h(\mathbf{x})$  is known for a given sample, the optimization problem can be written as follows by taking  $\|h(\mathbf{x})\| = c$ :

$$\begin{aligned} \min \quad & \frac{c + a \cos \theta}{\sqrt{(c + a \cos \theta)^2 + (a \sin \theta)^2}} \\ \text{s.t.} \quad & a \leq M \frac{V_U}{V_S} \end{aligned} \quad (29)$$

We now consider the Lagrange function of the optimization problem:

$$\mathcal{L}(x, \theta, \lambda) = \frac{c + a \cos \theta}{\sqrt{(c + a \cos \theta)^2 + (a \sin \theta)^2}} - \lambda(a - M \frac{V_U}{V_S}) \quad (30)$$

Taking the derivative of  $\mathcal{L}$  with respect to  $a$  and  $\theta$  and setting them to zero, we have

$$\frac{\partial}{\partial a} \mathcal{L} = \frac{1}{T^2} \left( T \cos \theta - \frac{1}{T} (c \cos \theta + 2a) \times (c + a \cos \theta) \right) - \lambda = 0 \quad (31)$$

<sup>2</sup> $\|\cdot\|$  in the following content denotes the  $\ell_2$ -norm unless otherwise specified.

810 and

$$811 \quad \frac{\partial}{\partial \theta} \mathcal{L} = \frac{1}{T^2} \left( -a \sin \theta \cdot T + \frac{1}{T} (c^2 a \sin \theta + ca^2 \sin \theta \cos \theta) \right) = 0 \quad (32)$$

812 where  $T = \sqrt{(c + a \cos \theta)^2 + (a \sin \theta)^2}$ . Solving the above equations, we have

$$813 \quad \cos \theta = 0 \quad \text{or} \quad a = 0 \quad (33)$$

814 where  $a = 0$  reaches the maximum and  $\cos \theta = 0$  is the minimum. Therefore, the lower bound of

815  $\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}}))$  is

$$816 \quad \cos(h(\mathbf{x}), h(\tilde{\mathbf{x}})) \geq \frac{c}{\sqrt{c^2 + (M \frac{V_U}{V_S})^2}} = \frac{\|h(\mathbf{x})\|}{\sqrt{\|h(\mathbf{x})\|^2 + (MV_U/V_S)^2}} \quad (34)$$

817  $\square$

818 **Corollary 1.** Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a bounded attribution function, and  $\eta \stackrel{U}{\sim} \mathcal{B}(\mathbf{x}; r)$ . Let  $h$  be the

819 smoothed version of  $g$  as defined in (2).

820 (i) Given a predefined threshold  $T \in [0, 1]$ , then for all  $\|\delta\|_2 \leq \epsilon$ , we have  $\cos(h(\mathbf{x}), h(\mathbf{x} + \delta)) \geq$

821  $T$ , where

$$822 \quad \epsilon = 2r \sqrt{1 - I_Z^{-1} \left( \frac{d+1}{2}, \frac{1}{2} \right)}. \quad (8)$$

823 (ii) Given a predefined threshold  $T \in [0, 1]$  and the maximum perturbation size  $\epsilon \geq 0$ , the

824 smoothed attribution satisfies  $\cos(h(\mathbf{x}), h(\tilde{\mathbf{x}})) \geq T$  for all  $\tilde{\mathbf{x}} \in \{\mathbf{x} + \delta \mid \|\delta\|_2 \leq \epsilon\}$  when

825  $r \geq R$ , where

$$826 \quad R = \frac{\epsilon}{2} \left( 1 - I_Z^{-1} \left( \frac{d+1}{2}, \frac{1}{2} \right) \right)^{-\frac{1}{2}}. \quad (9)$$

827  $I_z^{-1}(a, b)$  is the inverse of the regularized incomplete beta function, and  $Z$  is defined as

$$828 \quad Z = 1 - \frac{\|h(\mathbf{x})\|_2}{2M} \left( \frac{1}{T^2} - 1 \right) \quad (10)$$

829 *Proof.* Corollary 1 can be obtained by fixing  $T$  and taking  $r$  as unknown, and fixing  $T$  and taking  $\epsilon$

830 as unknown, respectively. We can first derive that

$$831 \quad I_{(2rh-h^2)/r^2} \left( \frac{d+1}{2}, \frac{1}{2} \right) = 1 - \frac{\|h(x)\|_2}{2M} \sqrt{\frac{1}{T^2} - 1} = Z \quad (35)$$

832 Using the inverse of the regularized incomplete beta function, i.e.,  $x = I_y^{-1}(a, b)$ , and  $h = r - \epsilon/2$ ,

833 we have

$$834 \quad I_Z^{-1} \left( \frac{d+1}{2}, \frac{1}{2} \right) = (2rh - h^2)/r^2 = 1 - \frac{\epsilon^2}{4r^2} \quad (36)$$

835 The results in Corollary can then be solved accordingly.  $\square$

## 836 B IMPLEMENTATION DETAILS

837 In the experiments, we implemented the  $\ell_2$  attribution attack adapted from Ghorbani et al. (2019).

838 The attack uses top- $k$  intersection version as the loss function. Following previous works, we choose

839  $k = 100$  for MNIST and  $k = 1000$  for CIFAR-10. The number of iterations in PGD-like attack is

840 200, and the step size is 0.1. As mentioned in the main content, we do not implement the attack on

841 ImageNet since the attribution attacks are not scalable to large size images. In the following parts of

842 this section, we provide more details of evaluations in the experiments.



## B.1 ATTRIBUTION METHODS

We used saliency maps (SM) and integrated gradients (IG) in the evaluation sections. These two methods are defined as follows:

- Saliency maps:  $\text{SM}(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$ .
- Integrated gradients:  $\text{IG}(\mathbf{x}) = (\mathbf{x} - \mathbf{x}') \times \int_{\alpha=0}^1 \frac{\partial f(\mathbf{x}' + \alpha(\mathbf{x} - \mathbf{x}'))}{\partial \mathbf{x}} d\alpha$ .

The SmoothSM and SmoothIG are the smoothed versions of SM and IG, respectively.

## B.2 EVALUATION METRICS

Given original attribution  $g(\mathbf{x})$  and perturbed attribution  $g(\tilde{\mathbf{x}})$ , we use top-k intersection, Kendall’s rank correlation (Ghorbani et al., 2019) and cosine similarity (Wang & Kong, 2022) to evaluate their differences.

- Top-k intersection measures the proportion of  $k$  largest features that overlap between  $g(\mathbf{x})$  and  $g(\tilde{\mathbf{x}})$ .
- Kendall’s rank correlation measures the proportion of pairs of features that have the same order in  $g(\mathbf{x})$  and  $g(\tilde{\mathbf{x}})$ :  $\frac{2}{d(d-1)} \sum_{i=1}^d \sum_{j=i+1}^d \mathbf{1}_{\{g(\mathbf{x})_i > g(\mathbf{x})_j\}} \mathbf{1}_{\{g(\tilde{\mathbf{x}})_i > g(\tilde{\mathbf{x}})_j\}}$ .
- Cosine similarity measures the cosine of the angle between  $g(\mathbf{x})$  and  $g(\tilde{\mathbf{x}})$ :  $\frac{g(\mathbf{x})^\top g(\tilde{\mathbf{x}})}{\|g(\mathbf{x})\| \|g(\tilde{\mathbf{x}})\|}$ .

## B.3 BASELINE METHODS

We compare with the following adversarial and attributional robust models:

### IG-NORM (Chen et al., 2019)

$$\text{CE}(f(\mathbf{x}), y) + \lambda \max_{\tilde{\mathbf{x}} \in \mathcal{B}_\varepsilon(\mathbf{x})} \|\text{IG}(\mathbf{x}, \tilde{\mathbf{x}})\|_1 \quad (37)$$

### TRADES (Zhang et al., 2019)

$$\text{CE}(f(\tilde{\mathbf{x}}), y) + \beta \text{KL}(f(\mathbf{x}) \| f(\tilde{\mathbf{x}})) \quad (38)$$

### IGR (Wang & Kong, 2022)

$$\text{CE}(f(\tilde{\mathbf{x}}), y) + \beta \text{KL}(f(\mathbf{x}) \| f(\tilde{\mathbf{x}})) + \lambda (1 - \cos(\text{IG}(\mathbf{x}), \text{IG}(\tilde{\mathbf{x}}))) \quad (39)$$

Here CE denotes the cross-entropy loss and KL denotes the Kullback-Leibler divergence.

## C ADDITIONAL EXPERIMENTS

### C.1 TEST ON MONTE CARLO ESTIMATION

Note that the bound given by Theorem 1 is deterministic. In this section, we provide a probabilistic bound for the attribution robustness. Specifically, we want to find the value of  $t$  such that  $\Pr(T \leq t) = 1 - \alpha$ , where  $T$  is defined in Eqn. (6) and  $\alpha$  is the significance level. Recall that  $T$  is defined as follows:

$$T = \frac{\|h(\mathbf{x})\|_2}{\sqrt{\|h(\mathbf{x})\|_2^2 + c}} \quad (40)$$

where  $c = M^2 V_U^2 / V_S^2$ . If we denote that  $Q = \|h(\mathbf{x})\|_2$ , then we have

$$\Pr(T \leq t) = \Pr\left(\frac{Q}{\sqrt{Q^2 + c}} \leq t\right) = \Pr\left(Q^2 \leq \frac{ct^2}{1 - t^2}\right) \quad (41)$$

Table 4: Evaluation of center smoothing on attributions

$\epsilon_1$	0.1	0.2	0.3	0.4	0.5
SmoothSM	1.207	1.729	1.843	1.907	1.998

Note that we used Monte Carlo Integration to calculate the integral in  $h(\mathbf{x})$ , which estimates  $h(\mathbf{x})$  by sampling  $\boldsymbol{\eta}$  from  $\mathcal{B}$ , *i.e.*,

$$\hat{h}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N g(\mathbf{x} + \boldsymbol{\eta}_i), \quad \boldsymbol{\eta}_i \sim \mathcal{B}. \quad (42)$$

Note that  $\hat{h}(\mathbf{x})$  is an unbiased estimator of  $h(\mathbf{x})$ , *i.e.*  $\mathbb{E}[\hat{h}(\mathbf{x})] = h(\mathbf{x})$ . The estimator almost surely converges to  $h(\mathbf{x})$  as  $N \rightarrow \infty$ , *i.e.*  $\lim_{N \rightarrow \infty} \hat{h}(\mathbf{x}) = h(\mathbf{x})$  almost surely. By the Central Limit Theorem, the estimator  $\hat{h}(\mathbf{x})$  has the following asymptotic distribution,

$$\hat{h}(\mathbf{x}) \stackrel{a.s.}{\sim} \mathcal{N}(h(\mathbf{x}), D), \quad (43)$$

which the covariance matrix  $D = \text{diag}(\sigma_{ii}^2/N)$  can be estimated by the empirical variances of  $g(\mathbf{x} + \boldsymbol{\eta}_i)$ . Thus, the quadratic form  $Q^2 = \|\hat{h}(\mathbf{x})\|_2^2$  can be seen as generalized chi-square distributed. We can derive the cumulative distribution function of Monte Carlo estimator  $T_{MC}$  at  $t$  as the cumulative distribution function of the generalized chi-square distribution at  $\frac{ct^2}{1-t^2}$ , *i.e.*,

$$Pr(T_{MC} \leq t) = F\left(\frac{ct^2}{1-t^2}\right), \quad (44)$$

where  $F$  is the cumulative distribution function of the generalized chi-square distribution constructed from the quadratic form of Gaussian random variable with mean  $h(\mathbf{x})$  and covariance  $D$  (Davies, 1980; Das & Geisler, 2021). In this work, we use the R package `CompQuadForm` (Duchesne & De Micheaux, 2010) to compute the cumulative distribution function. For any fixed image sample  $\mathbf{x}$ , we can validate  $t_2 - t_1$  is close to 0 when  $Pr(t_1 \leq T_{MC} \leq t_2) = 1 - \alpha$  by solving the following equation. For small  $\alpha = 0.01$  and the number of samples  $N = 100,000$ , we found that the values of  $t_2 - t_1$  are at scale of  $10^{-4}$  in MNIST and CIFAR-10, and  $10^{-3}$  in ImageNet calculated by choosing 10,000 samples from each dataset. This validates the error from Monte Carlo integral is minute and that the probabilistic bound is close to the deterministic bound.

$$F\left(\frac{ct_2^2}{1-t_2^2}\right) = 1 - \alpha/2 \quad \text{and} \quad F\left(\frac{ct_1^2}{1-t_1^2}\right) = \alpha/2. \quad (45)$$

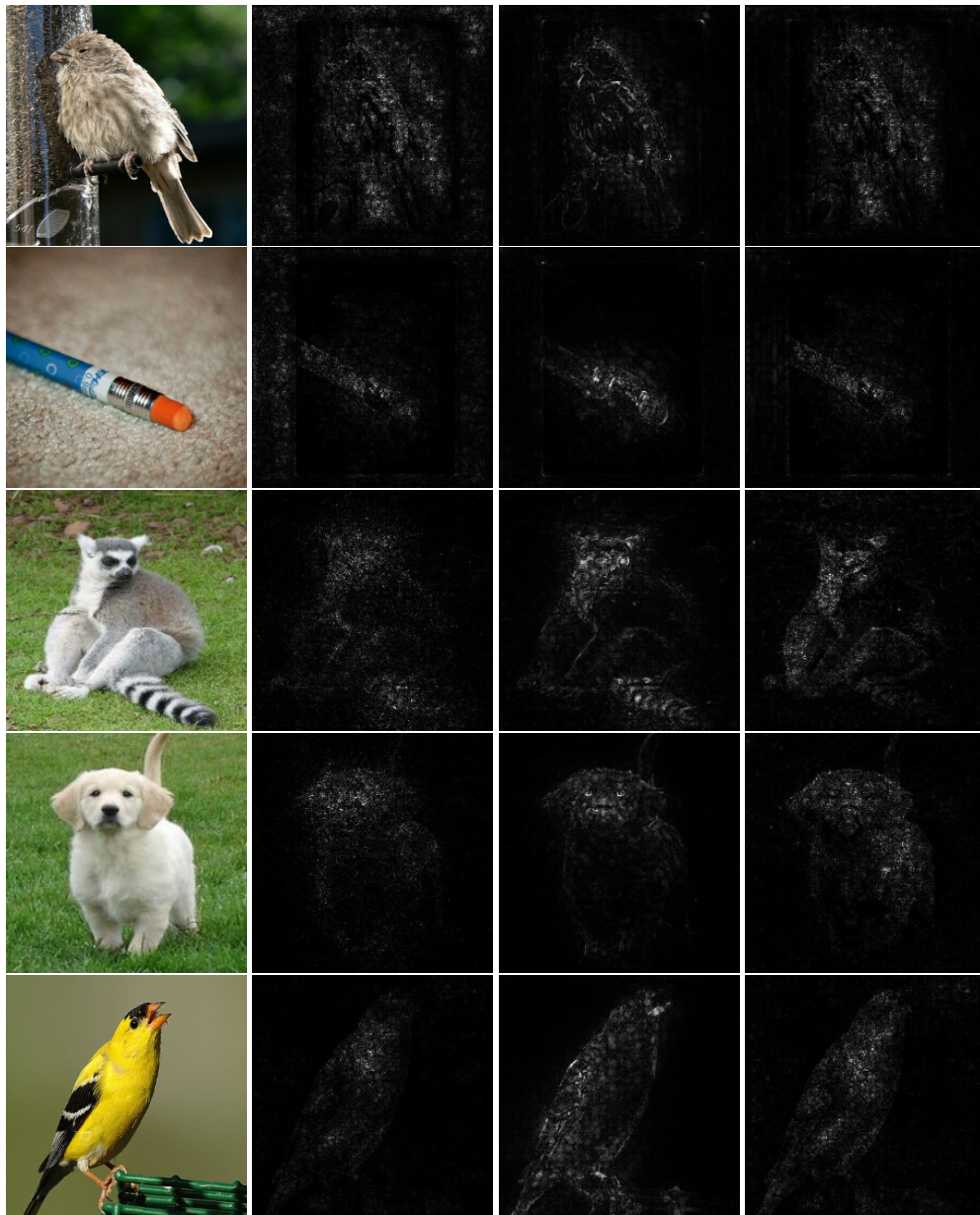
## C.2 ADDITIONAL VISUALIZATION OF THE UNIFORMLY SMOOTHED ATTRIBUTIONS

In Figure 1 (left), we have shown that the uniformly smoothed attributions have a comparable quality as the original attributions. Here more examples are provided in Figure 3 to illustrate the quality of the uniformly smoothed attributions.

## C.3 EVALUATION OF CENTER SMOOTHING (KUMAR & GOLDSTEIN, 2021) ON ATTRIBUTIONS

To compare the performance with center smoothing (Kumar & Goldstein, 2021), we also implemented the same method to evaluate the certification of attributions. Specifically, we compute the bound for SmoothSM on IG-NORM using MNIST, and follow the same setting by choosing  $h = 1$  and  $\epsilon_1 = 0.1, 0.2, \dots, 0.5$ . Directly using the cosine similarity on the method is not applicable since cosine similarity does not satisfy the triangle inequality. Following the relaxation method in Sec.4 of Kumar & Goldstein (2021), a multiplier  $\gamma = 2$  is added. Besides, we use  $1 - \cos \theta$  to reflect the distance metric instead of the similarity metric. The results are shown in the Table 4. It can be observed that the upper bound for  $1 - \cos \theta$  is greater than 1 for all the choices of  $\epsilon$ , which is trivially valid for the trigonometric function since we only consider  $\cos \theta \in [0, 1]$ . Thus, the upper bound provided in the aforementioned work can be too loose on our setting.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010



(a) Original image (b) IG (c) Gaussian smoothed (d) Uniformly smoothed

1011  
1012  
1013  
1014

Figure 3: Additional visualization of the attribution maps of the (a) original image, (b) IG, (c) Gaussian smoothed IG, and (d) uniformly smoothed IG.

1015  
1016

#### 1017 C.4 EVALUATION OF ALTERNATIVE FORMULATIONS

1018  
1019  
1020  
1021

In Section 4.3, we introduced two alternative formulations of the proposed method that can be applied in specific scenarios. In this section, we provide additional information to report the experiments on these two formulations.

1022  
1023  
1024  
1025

In Tables 5 to 7, which correspond to MNIST, CIFAR-10 and ImageNet, respectively, we report the computed values of the maximum allowable perturbation size. Under the size constraint, no examples can be found by the attacks against uniformly smoothed IG of a certain radius such that the cosine similarity between clean and perturbed attributions exceeds the given threshold ( $T = 0.8$  and  $T = 0.9$ ). The results are consistent with our theory. For larger radius smoothing, the maximum

Table 5: Maximum allowable perturbation size for different threshold ( $T = 0.8$  and  $T = 0.9$ ) under various choices of  $\ell_2$  smoothing radii  $r$  evaluated on MNIST.

$T = 0.9$	$\ell_2$ radius ( $r$ )	0.5	1.0	1.5	2.0	2.5	3.0	3.5
	Standard	0.0389	0.0951	0.1550	0.2164	0.2783	0.3404	0.4029
	IG-NORM	0.0394	0.0957	0.1557	0.2170	0.2790	0.3420	0.4067
	IGR	0.0390	0.0952	0.1552	0.2174	0.2818	0.3477	0.4163
$T = 0.8$	$\ell_2$ radius ( $r$ )	0.5	1.0	1.5	2.0	2.5	3.0	3.5
	Standard	0.0447	0.1051	0.1691	0.2345	0.3004	0.3664	0.4329
	IG-NORM	0.0448	0.1052	0.1692	0.2354	0.3037	0.3733	0.4456
	IGR	0.0452	0.1057	0.1697	0.2350	0.3010	0.3680	0.4365

Table 6: Maximum allowable perturbation size for different threshold ( $T = 0.8$  and  $T = 0.9$ ) under various choices of  $\ell_2$  smoothing radii  $r$  evaluated on CIFAR-10.

$T = 0.9$	$\ell_2$ radius ( $r$ )	0.5	1.0	1.5	2.0	2.5	3.0	3.5
	Standard	0.0086	0.0469	0.0885	0.1322	0.1773	0.2222	0.2683
	IG-NORM	0.0323	0.0705	0.1104	0.1510	0.1923	0.2337	0.2749
	IGR	0.0167	0.0545	0.1032	0.1586	0.2150	0.2588	0.2805
$T = 0.8$	$\ell_2$ radius ( $r$ )	0.5	1.0	1.5	2.0	2.5	3.0	3.5
	Standard	0.0128	0.0522	0.0951	0.1402	0.1866	0.2330	0.2868
	IG-NORM	0.0343	0.0742	0.1157	0.1580	0.2009	0.2439	0.2867
	IGR	0.0237	0.0693	0.1258	0.1861	0.2546	0.3090	0.3559

Table 7: Maximum allowable perturbation size for different threshold ( $T = 0.8$  and  $T = 0.9$ ) under various choices of  $\ell_2$  smoothing radii  $r$  evaluated on ImageNet.

$\ell_2$ radius ( $r$ )	0.5	1.0	1.5	2.0	2.5	3.0	3.5
$T = 0.9$	0.0046	0.0100	0.0152	0.0295	0.0494	0.0628	0.0768
$T = 0.8$	0.0058	0.0127	0.0196	0.0369	0.0618	0.0820	0.1040

Table 8: Empirical cosine similarity between original and perturbed smoothed attributions under various choices of  $\ell_2$  smoothing radius  $r$ , and the perturbation size computed in Table 5 ( $T = 0.8$ ).

$r$	0.5	1.0	1.5	2.0	2.5	3.0	3.5
Standard	0.8636	0.8522	0.8347	0.8127	0.8477	0.8603	0.8310
IG-NORM	0.8308	0.8181	0.8504	0.8728	0.8502	0.8193	0.8199
IGR	0.8231	0.8800	0.8720	0.8603	0.8362	0.8135	0.8567

Table 9: Minimum smoothing radius requires to achieve the threshold ( $T = 0.8$  and  $T = 0.9$ ) under various choices of  $\ell_2$  perturbation size  $\epsilon$ . IG-NORM and IGR are omitted since they are not scalable to ImageNet.

		MNIST		CIFAR-10		ImageNet	
perturbation size ( $\epsilon$ )		0.5	1.0	0.5	1.0	0.5	1.0
$T = 0.9$	Standard	5.1902	5.8752	5.9752	7.9504	74.6272	149.2544
	IG-NORM	5.1189	5.7699	5.6860	7.3720	/	/
	IGR	5.0265	5.6623	5.2895	6.5790	/	/
$T = 0.8$	Standard	3.8927	4.4064	5.7082	7.4164	48.2095	96.4190
	IG-NORM	3.8392	4.3274	5.4875	6.9750	/	/
	IGR	3.7699	4.2468	5.0287	6.0573	/	/

1080 allowable perturbation size is also larger. When the threshold requirement is stricter, the maximum  
1081 allowable perturbation size is smaller, which suggests weaker attacks are allowed. The method is  
1082 also scalable to ImageNet, which takes around 15 seconds to compute for each sample. Moreover,  
1083 we also applied attribution attacks using the same radius and maximum perturbation size  $\epsilon$ , computed  
1084 using Eqn. (8). Similar to the experiments in Section 5, we performed 20 attacks on each sample.  
1085 We found that out of the total 200,000 attacked samples, the cosine similarities between clean and  
1086 perturbed attributions were higher than the given threshold, suggesting that the computed bound is  
1087 valid (see Table 8).

1088 We also evaluate the third formulation that the minimum radius of smoothing required such that,  
1089 within the given perturbation sizes, the cosine similarity between original and perturbed smoothed  
1090 attributions is larger than the given threshold. In Table 9, the computed minimum radius of smooth-  
1091 ing is reported. Similarly, we observe that the minimum radius of smoothing is larger when the  
1092 threshold requirement is stricter, and when the attack is stronger. This is also consistent with our  
1093 theory. We also notice that the radius for ImageNet is extremely large, which indicates that ImageNet  
1094 is difficult to defend under such strict threshold requirements.

1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133