

LOCALIZED RANDOMIZED SMOOTHING FOR COLLECTIVE ROBUSTNESS CERTIFICATION

Jan Schuchardt^{1*}, Tom Wollschläger^{1*}, Aleksandar Bojchevski², Stephan Günnemann¹
 {j.schuchardt, t.wollschlaeger, s.guennemann}@tum.de
 {bojchevski}@cispa.de

¹Technical University of Munich

²CISPA Helmholtz Center for Information Security

ABSTRACT

Models for image segmentation, node classification and many other tasks map a single input to multiple labels. By perturbing this single shared input (e.g. the image) an adversary can manipulate several predictions (e.g. misclassify several pixels). Collective robustness certification is the task of provably bounding the number of robust predictions under this threat model. The only dedicated method that goes beyond certifying each output independently is limited to *strictly local* models, where each prediction is associated with a small receptive field. We propose a more general collective robustness certificate for all types of models. We further show that this approach is beneficial for the larger class of *softly local* models, where each output is dependent on the entire input but assigns different levels of importance to different input regions (e.g. based on their proximity in the image). The certificate is based on our novel localized randomized smoothing approach, where the random perturbation strength for different input regions is proportional to their importance for the outputs. Localized smoothing Pareto-dominates existing certificates on both image segmentation and node classification tasks, simultaneously offering higher accuracy and stronger certificates.

1 INTRODUCTION

There is a wide range of tasks that require models making multiple predictions based on a single input. For example, semantic segmentation requires assigning a label to each pixel in an image. When deploying such *multi-output* classifiers in practice, their robustness should be a key concern. After all – just like simple classifiers (Szegedy et al., 2014) – they can fall victim to adversarial attacks (Xie et al., 2017; Zügner & Günnemann, 2019; Belinkov & Bisk, 2018). Even without an adversary, random noise or measuring errors can cause predictions to unexpectedly change.

We propose a novel method providing provable guarantees on *how many* predictions can be changed by an adversary. As all outputs operate on the same input, they have to be attacked simultaneously by choosing a single perturbed input, which can be more challenging for an adversary than attacking them independently. We must account for this to obtain a proper *collective robustness certificate*.

The only dedicated collective certificate that goes beyond certifying each output independently (Schuchardt et al., 2021) is only beneficial for models we call *strictly local*, where each output depends on a small, pre-defined subset of the input. Multi-output classifiers, however, are often only *softly local*. While all their predictions are in principle dependent on the entire input, each output may assign different importance to different subsets. For example, convolutional networks for image segmentation can have small effective receptive fields (Luo et al., 2016; Liu et al., 2018), i.e. primarily use a small region of the image in labeling each pixel. Many models for node classification are based on the homophily assumption that connected nodes are mostly of the same class. Thus, they primarily use features from neighboring nodes. Transformers, which can in principle attend to arbitrary parts of the input, may in practice learn “sparse” attention maps, with the prediction for each token being mostly determined by a few (not necessarily nearby) tokens (Shi et al., 2021).

*equal contribution

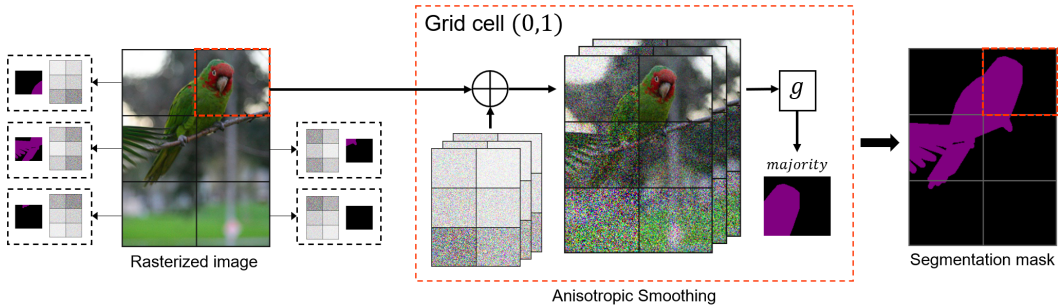


Figure 1: Localized randomized smoothing applied to semantic segmentation. We assume that the most relevant information for labeling a pixel is contained in other nearby pixels. We partition the input image into multiple grid cells. For each grid cell, we sample noisy images from a *different* anisotropic distribution that applies more noise to far-away, less relevant cells. Segmenting all noisy images, cropping the result and computing the majority vote yields a local segmentation mask. These per-cell segmentation masks can then be combined into a complete segmentation mask.

Softly local models pose a *budget allocation problem* for an adversary that tries to simultaneously manipulate multiple predictions by crafting a single perturbed input. When each output is primarily focused on a different part of the input, the attacker has to distribute their limited adversarial budget and may be unable to attack all predictions at once.

We propose *localized randomized smoothing*, a novel method for the collective robustness certification of softly local models that exploits this budget allocation problem. It is an extension of randomized smoothing (Lécuyer et al., 2019; Li et al., 2019; Cohen et al., 2019), a versatile black-box certification method which is based on constructing a smoothed classifier that returns the expected prediction of a model under random perturbations of its input (more details in § 2). Randomized smoothing is typically applied to single-output models with *isotropic* Gaussian noise. In localized smoothing however, we smooth each output (or set of outputs) of a multi-output classifier using a *different* distribution that is *anisotropic*. This is illustrated in Fig. 1, where the predicted segmentation masks for each grid cell are smoothed using a different distribution. For instance, the distribution for segmenting the top-right cell applies less noise to the top-right cell. The smoothing distribution for segmenting the bottom-left cell applies significantly more noise to the top-right cell.

Given a specific output of a softly local model, using a low noise level for the most relevant parts of the input lets us preserve a high prediction quality. Less relevant parts can be smoothed with a higher noise level to guarantee more robustness. The resulting certificates (one per output) explicitly quantify how robust each prediction is to perturbations of which part of the input. This information about the smoothed model’s locality can then be used to combine the per-prediction certificates into a stronger collective certificate that accounts for the adversary’s budget allocation problem.¹

Our core contributions are:

- *Localized randomized smoothing*, a novel smoothing scheme for multi-output classifiers.
- An efficient anisotropic randomized smoothing certificate for discrete data.
- A collective certificate based on localized randomized smoothing.

2 BACKGROUND AND RELATED WORK

Randomized smoothing. Randomized smoothing is a certification technique that can be used for various threat models and tasks. For the sake of exposition, let us discuss a certificate for l_2 perturbations (Cohen et al., 2019). Assume we have a D -dimensional input space \mathbb{R}^D , label set \mathbb{Y} and classifier $g : \mathbb{R}^D \rightarrow \mathbb{Y}$. We can use isotropic Gaussian noise to construct the *smoothed classifier* $f = \operatorname{argmax}_{y \in \mathbb{Y}} \Pr_{z \sim \mathcal{N}(x, \sigma)} [g(z) = y]$ that returns the most likely prediction of *base classifier* g under the input distribution². Given an input $x \in \mathbb{R}^D$ and smoothed prediction $y = f(x)$, we can then easily determine whether y is robust to all l_2 perturbations of magnitude ϵ , i.e. whether $\forall x' : \|x' - x\|_2 \leq \epsilon : f(x') = y$. Let $q = \Pr_{z \sim \mathcal{N}(x, \sigma)} [g(z) = y]$ be the probability of predicting

¹An implementation will be made available at <https://www.cs.cit.tum.de/daml/localized-smoothing>.

²In practice, all probabilities have to be estimated using Monte Carlo sampling (see discussion in § G).

label y . The prediction is certifiably robust if $\epsilon < \sigma\Phi^{-1}(q)$ (Cohen et al., 2019). This result shows a trade-off inherent to randomized smoothing: Increasing the noise level (σ) may strengthen the certificate, but could also lower the accuracy of f or reduce q and thus weaken the certificate.

White-box certificates for multi-output classifiers. There are multiple recent methods for certifying the robustness of multi-output models by analyzing their specific architecture and weights (for example, see (Tran et al., 2021; Zügner & Günnemann, 2019; Bojchevski & Günnemann, 2019; Zügner & Günnemann, 2020; Ko et al., 2019; Ryou et al., 2021; Shi et al., 2020; Bonaert et al., 2021)). They are however not designed to certify collective robustness, i.e. determine whether multiple outputs can be simultaneously attacked using a single perturbed input. They can only determine independently for each prediction whether or not it can be attacked.

Black-box certificates for multi-output classifiers. Most directly related to our work is the aforementioned certificate of Schuchardt et al. (2021), which is only beneficial for strictly local models (i.e. models where each output has a small receptive field). In § I we show that, for randomly smoothed models, their certificate is a special case of ours. SegCertify (Fischer et al., 2021) is a collective certificate for segmentation. This method certifies each output independently using isotropic smoothing (ignoring the budget allocation problem) and uses Holm correction (Holm, 1979) to obtain tighter Monte Carlo estimates. It then counts the number of certifiably robust predictions and tests whether it equals the number of predictions. In § H we demonstrate that our method can always provide guarantees that are at least as strong. Another method that can in principle be used to certify collective robustness is center smoothing (Kumar & Goldstein, 2021). It bounds the change of a vector-valued function w.r.t to a distance function. Using the l_0 pseudo-norm, it can bound how many predictions can be simultaneously changed. More recently, Chen et al. (2022) proposed a collective certificate for bagging classifiers. Different from our work, they consider poisoning (train-time) instead of evasion (test-time) attacks. Yatsura et al. (2022) prove robustness for segmentation, but consider patch-based instead of ℓ_p -norm attacks and certify each prediction independently.

Anisotropic randomized smoothing. While only designed for single-output classifiers, two recent certificates for anisotropic Gaussian and uniform smoothing (Fischer et al., 2020; Eiras et al., 2022) can be used as a component of our collective certification approach: They can serve as per-prediction certificates, which we can then combine into our stronger collective certificate (more details in § 3.2).

3 PRELIMINARIES

3.1 COLLECTIVE THREAT MODEL

We assume a multi-output classifier $f : \mathbb{X}^{D_{\text{in}}} \rightarrow \mathbb{Y}^{D_{\text{out}}}$, that maps D_{in} -dimensional inputs to D_{out} labels from label set \mathbb{Y} . We further assume that this classifier f is the result of randomly smoothing each output of a base classifier g . Given this multi-output classifier f , an input $\mathbf{x} \in \mathbb{X}^{D_{\text{in}}}$ and the corresponding predictions $\mathbf{y} = f(\mathbf{x})$, the objective of the adversary is to cause as many predictions from a set of targeted indices $\mathbb{T} \subseteq \{1, \dots, D_{\text{out}}\}$ to change. That is, their objective is $\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(\mathbf{x}') = y_n]$, where \mathbb{I} is the indicator function and $\mathbb{B}_{\mathbf{x}} \subseteq \mathbb{X}^{D_{\text{in}}}$ is the perturbation model. As is common in robustness certification, we assume a ℓ_p -norm perturbation model, i.e. $\mathbb{B}_{\mathbf{x}} = \{\mathbf{x}' \in \mathbb{X}^{D_{\text{in}}} \mid \|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon\}$ with $p, \epsilon \geq 0$. Importantly, note that the minimization operator is outside the sum, meaning the predictions have to be attacked using a single input.

3.2 A RECIPE FOR COLLECTIVE CERTIFICATES

Before discussing localized randomized smoothing, we show how to combine arbitrary per-prediction certificates into a collective certificate, a procedure that underlies both our method and that of Schuchardt et al. (2021) and Fischer et al. (2021). The first step is to apply an arbitrary certification procedure to each prediction $y_1, \dots, y_{D_{\text{out}}}$ in order to obtain per-prediction *base certificates*.

Definition 3.1 (Base certificates). A base certificate for a prediction $y_n = f_n(\mathbf{x})$ is a set $\mathbb{H}^{(n)} \subseteq \mathbb{X}^{D_{\text{in}}}$ of perturbed inputs s.t. $\forall \mathbf{x}' \in \mathbb{H}^{(n)} : f_n(\mathbf{x}') = y_n$.

Using these base certificates, one can derive two bounds on the adversary’s objective:

$$\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(\mathbf{x}') = y_n] \stackrel{(1.1)}{\geq} \min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[\mathbf{x}' \in \mathbb{H}^{(n)}] \stackrel{(1.2)}{\geq} \sum_{n \in \mathbb{T}} \min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \mathbb{I}[\mathbf{x}' \in \mathbb{H}^{(n)}]. \quad (1)$$

Eq. 1.1 follows from Theorem 3.1 (if a prediction is certifiably robust to \mathbf{x}' , then $f_n(\mathbf{x}') = y_n$), while Eq. 1.2 results from moving the min operator inside the summation.

Eq. 1.2 is the *naïve collective certificate*: It iterates over the predictions and counts how many are certifiably robust to perturbation model $\mathbb{B}_{\mathbf{x}}$. Each summand involves a separate minimization problem. Thus, the certificate neglects that the adversary has to choose a single perturbed input to attack all outputs. SegCertify (Fischer et al., 2021) applies this to isotropic Gaussian smoothing.

While Eq. 1.1 is seemingly tighter than the naïve collective certificate, it may lead to identical results. For example, let us consider the most common case where the base certificates guarantee robustness within an l_p ball, i.e. $\mathbb{H}^{(n)} = \{\mathbf{x}'' \mid \|\mathbf{x}'' - \mathbf{x}\|_p \leq r^{(n)}\}$ with certified radii $r^{(n)}$. Then, the optimal solution to both Eq. 1.1 and Eq. 1.2 is to choose an arbitrary \mathbf{x}' with $\|\mathbf{x}' - \mathbf{x}\| = \epsilon$:

$$\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I} \left[\mathbf{x}' \in \mathbb{H}^{(n)} \right] = \sum_{n \in \mathbb{T}} \mathbb{I} \left[\epsilon < r^{(n)} \right] = \sum_{n \in \mathbb{T}} \min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \mathbb{I} \left[\mathbf{x}' \in \mathbb{H}^{(n)} \right].$$

The main contribution of Schuchardt et al. (2021) is to notice that, by exploiting strict locality (i.e. the outputs having small receptive fields), one can augment certificate Eq. 1.1 to make it tighter than the naïve collective certificate from Eq. 1.2. One must simply mask out all perturbations falling outside a given receptive field when evaluating the corresponding base certificate:

$$\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I} \left[\left(\boldsymbol{\psi}^{(n)} \odot \mathbf{x}' + (1 - \boldsymbol{\psi}^{(n)}) \odot \mathbf{x} \right) \in \mathbb{H}^{(n)} \right].$$

Here, $\boldsymbol{\psi}^{(n)} \in \{0, 1\}^{D_{\text{in}}}$ encodes the receptive field of f_n and \odot is the elementwise product. If two outputs f_n and f_m have disjoint receptive fields (i.e. $\boldsymbol{\psi}^{(n)T} \boldsymbol{\psi}^{(m)} = 0$), then the adversary has to split up their limited adversarial budget and may be unable to attack both at once.

4 LOCALIZED RANDOMIZED SMOOTHING

The core idea behind localized smoothing is that, rather than improving upon the naïve collective certificate by using external knowledge about *strict* locality, we can use anisotropic randomized smoothing to obtain base certificates that directly encode *soft* locality. Here, we explain our approach in a domain-independent manner before turning to specific distributions and data-types in § 5.

In localized randomized smoothing, we associate base classifier outputs $g_1, \dots, g_{D_{\text{out}}}$ with distinct anisotropic smoothing distributions $\Psi_{\mathbf{x}}^{(1)}, \dots, \Psi_{\mathbf{x}}^{(D_{\text{out}})}$ that depend on input \mathbf{x} . For example, they could be Gaussian distributions with mean \mathbf{x} and distinct covariance matrices – like in Fig. 1, where we use a different distribution for each grid cell. We use these distributions to construct the smoothed classifier f , where each output $f_n(\mathbf{x})$ is the result of randomly smoothing $g_n(Z)$ with $\Psi_{\mathbf{x}}^{(n)}$.

To certify robustness for a vector of predictions $\mathbf{y} = f(\mathbf{x})$, we follow the procedure discussed in § 3.2, i.e. compute base certificates $\mathbb{H}^{(1)}, \dots, \mathbb{H}^{(D_{\text{out}})}$ and solve Eq. 1.1. We do not make any assumption about how the base certificates are computed. However, we require that they comply with a common interface, which will later allow us combine them via linear programming:

Definition 4.1 (Base certificate interface). A base certificate $\mathbb{H}^{(n)} \subseteq \mathbb{X}^{D_{\text{in}}}$ is compliant with our base certificate interface for l_p -norm perturbations if there is a $\mathbf{w} \in \mathbb{R}_+^{D_{\text{in}}}$ and $\eta^{(n)} \in \mathbb{R}_+$ such that

$$\mathbb{H}^{(n)} = \left\{ \mathbf{x}' \mid \sum_{d=1}^{D_{\text{in}}} w_d^{(n)} \cdot |x'_d - x_d|^p < \eta^{(n)} \right\}. \quad (2)$$

The weight $w_d^{(n)}$ quantifies how sensitive y_n is to perturbations of input dimension d . It will be smaller where the anisotropic smoothing distribution applies more noise. The radius $\eta^{(n)}$ quantifies the overall level of robustness. In § 5 we present different distributions and corresponding certificates that comply with this interface. Inserting Eq. 2 into Eq. 1.1 results in the collective certificate

$$\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I} \left[\sum_{d=1}^{D_{\text{in}}} w_d^{(n)} \cdot |x'_d - x_d|^p < \eta^{(n)} \right]. \quad (3)$$

Eq. 3 showcases why locally smoothed models admit a collective certificate that is stronger than naively certifying each output independently (i.e. Eq. 1.2). Because we use different distributions for different outputs, any two outputs $f^{(n)}$ and $f^{(m)}$ will have distinct certificate weights $\mathbf{w}^{(n)}$ and $\mathbf{w}^{(m)}$. If they are sensitive in different parts of the input, i.e. $\mathbf{w}^{(n)T}\mathbf{w}^{(m)}$ is small, then the adversary has to split up their limited adversarial budget and may be unable to attack both at once.

One particularly simple example is the case $\mathbf{w}^{(n)T}\mathbf{w}^{(m)} = 0$, where attacking predictions y_n and y_m requires allocating adversarial budget to two entirely disjoint sets of input dimensions. In § I we show that, with appropriately parameterized smoothing distributions, we can obtain base certificates with $\mathbf{w}^{(n)} = c \cdot \boldsymbol{\psi}^{(n)}$, with indicator vector $\boldsymbol{\psi}^{(n)}$ encoding the receptive field of output n . Hence, the collective guarantees from (Schuchardt et al., 2021) are a special case of our certificate.

4.1 COMPUTING THE COLLECTIVE CERTIFICATE

While Eq. 3 constitutes a valid certificate, it is not immediately clear how to evaluate it. However, we notice that the perturbation set $\mathbb{B}_{\mathbf{x}}$ imposes linear constraints on the elementwise differences $|x'_d - x_d|^p$, the values of the indicator functions are binary variables and that the base certificates inside the indicator functions are characterized by linear inequalities. We can thus reformulate Eq. 3 as a mixed-integer linear program (MILP), which leads us to our main result (proof in § D):

Theorem 4.2. *Given locally smoothed model f , input $\mathbf{x} \in \mathbb{X}^{(D_{\text{in}})}$, smoothed prediction $\mathbf{y} = f(\mathbf{x})$ and base certificates $\mathbb{H}^{(1)}, \dots, \mathbb{H}^{D_{\text{out}}}$ complying with interface Eq. 2, the number of simultaneously robust predictions $\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(\mathbf{x}') = y_n]$ is lower-bounded by*

$$\min_{\mathbf{b} \in \mathbb{R}_+^{D_{\text{in}}}, \mathbf{t} \in \{0,1\}^{D_{\text{out}}}} \sum_{n \in \mathbb{T}} t_n \quad (4)$$

$$\text{s.t. } \forall n : \mathbf{b}^T \mathbf{w}^{(n)} \geq (1 - t_n) \eta^{(n)}, \quad \text{sum}\{\mathbf{b}\} \leq \epsilon^p. \quad (5)$$

The vector \mathbf{b} models the allocation of adversarial budget (i.e. the elementwise differences $b_d = |x'_d - x_d|^p$). The vector \mathbf{t} serves the same role as the indicator functions from Eq. 3, i.e. it indicates which predictions are certifiably robust. Eq. 5 ensures that \mathbf{b} does not exceed the overall budget ϵ (i.e. $\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}$) and that t_n can only be set to 0 if $\mathbf{b}^T \mathbf{w}^{(n)} \geq \eta^{(n)}$, i.e. only when the base certificate cannot guarantee robustness for prediction y_n . This problem can be solved using any MILP solver. Its optimal value provably bounds the number of simultaneously robust predictions.

4.2 IMPROVING EFFICIENCY

Solving large MILPs is expensive. In § E we show that partitioning the outputs into N_{out} subsets sharing the same smoothing distribution and the inputs into N_{in} subsets sharing the same noise level (for example like in Fig. 1, where we partition the image into a 2×3 grid), as well as quantizing the base certificate parameters $\eta^{(n)}$ into N_{bins} bins, reduces the number of variables and constraints from $D_{\text{in}} + D_{\text{out}}$ and $D_{\text{out}} + 1$ to $N_{\text{in}} + N_{\text{out}} \cdot N_{\text{bins}}$ and $N_{\text{out}} \cdot N_{\text{bins}} + 1$, respectively. We can thus control the problem size independent of the data’s dimensionality. We further derive a linear relaxation of the MILP that can be efficiently solved while preserving the soundness of the certificate.

4.3 ACCURACY-ROBUSTNESS TRADEOFF

When discussing Eq. 3, we only explained why our collective certificate for *locally smoothed* models is better than a naïve combination of localized smoothing base certificates. However, this does not necessarily mean that our certificate is also stronger than naively certifying an *isotropically smoothed* model. This is why we focus on soft locality. With isotropic smoothing, high certified robustness requires using large noise levels, which degrade the model’s prediction quality. Localized smoothing, when applied to softly local models, can circumvent this issue. For each output, we can use low noise levels for the most important parts of the input to retain high prediction quality. Our LP-based collective certificate allows us to still provide strong collective robustness guarantees. We investigate this improved accuracy-robustness trade-off in our experimental evaluation (see § 7).

5 BASE CERTIFICATES

To apply our collective certificate in practice, we require smoothing distributions $\Psi_{\mathbf{x}}^{(n)}$ and corresponding per-prediction base certificates that comply with the interface from Theorem 3.1. As base certificates for l_2 and l_1 perturbations we can reformulate existing anisotropic Gaussian (Fischer et al., 2020; Kumar & Goldstein, 2021) and uniform (Kumar & Goldstein, 2021) smoothing certificates for single-output models: For $\Psi_{\mathbf{x}}^{(n)} = \mathcal{N}(\mathbf{x}, \text{diag}(\mathbf{s}^{(n)}))$ we have $w_d^{(n)} = 1/(s_d^{(n)})^2$ and $\eta^{(n)} = (\Phi^{-1}(q_{n,y_n}))^2$ with $q_{n,y_n} = \Pr_{\mathbf{z} \sim \Psi_{\mathbf{x}}^{(n)}} [g_n(\mathbf{z}) = y]$. For $\Psi_{\mathbf{x}}^{(n)} = \mathcal{U}(\mathbf{x}, \boldsymbol{\lambda}^{(n)})$ we have $w_d^{(n)} = 1/\lambda_d^{(n)}$ and $\eta^{(n)} = \Phi^{-1}(q_{n,y_n})$. We prove the correctness of these reformulations in § F.

For l_0 perturbations of binary data, we can use a distribution $\mathcal{F}(\mathbf{x}, \boldsymbol{\theta})$ that flips x_d with probability $\theta_d \in [0, 1]$, i.e. $\Pr[z_d \neq x_d] = \theta_d$ for $\mathbf{z} \sim \mathcal{F}(\mathbf{x}, \boldsymbol{\theta})$. Existing methods (e.g. (Lee et al., 2019)) can be used to derive per-prediction certificates for this distribution, but have exponential runtime in the number of unique values in $\boldsymbol{\theta}$. Thus, they are not suitable for localized smoothing, which uses different θ_d for different parts of the input. We therefore propose a novel, more efficient approach: *Variance-constrained certification*, which smooths the base classifier’s softmax scores instead of its predictions and then uses both their expected value and variance to certify robustness (proof in § F.3):

Theorem 5.1 (Variance-constrained certification). *Given a function $g : \mathbb{X} \rightarrow \Delta_{|\mathbb{Y}|}$ mapping from discrete set \mathbb{X} to scores from the $(|\mathbb{Y}| - 1)$ -dimensional probability simplex, let $f(\mathbf{x}) = \arg\max_{y \in \mathbb{Y}} \mathbb{E}_{\mathbf{z} \sim \Psi_{\mathbf{x}}} [g(\mathbf{z})_y]$ with smoothing distribution $\Psi_{\mathbf{x}}$ and probability mass function $\pi_{\mathbf{x}}(\mathbf{z}) = \Pr_{\tilde{\mathbf{z}} \sim \Psi_{\mathbf{x}}} [\tilde{\mathbf{z}} = \mathbf{z}]$. Given an input $\mathbf{x} \in \mathbb{X}$ and smoothed prediction $y = f(\mathbf{x})$, let $\mu = \mathbb{E}_{\mathbf{z} \sim \Psi_{\mathbf{x}}} [g(\mathbf{z})_y]$ and $\zeta = \mathbb{E}_{\mathbf{z} \sim \Psi_{\mathbf{x}}} [(g(\mathbf{z})_y - \nu)^2]$ with $\nu \in \mathbb{R}$. Assuming $\nu \leq \mu$, then $f(\mathbf{x}') = y$ if*

$$\sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})}{\pi_{\mathbf{x}}(\mathbf{z})} \cdot \pi_{\mathbf{x}'}(\mathbf{z}) < 1 + \frac{1}{\zeta - (\mu - \nu)^2} \left(\mu - \frac{1}{2} \right). \quad (6)$$

The l.h.s. of Eq. 6 is the expected ratio between the probability mass functions of the smoothing distributions for the perturbed ($\pi_{\mathbf{x}'}$) and unperturbed ($\pi_{\mathbf{x}}$) input.³ It is equal to 1 if both densities are the same, i.e. there is no adversarial perturbation, and greater than 1 otherwise. The r.h.s. of Eq. 6 depends on the expected softmax score μ , a variable $\nu \leq \mu$ and the expected squared difference ζ between μ and ν . For $\nu = \mu$ the parameter ζ is the variance of the softmax score. A higher expected value and a lower variance allow us to certify robustness for larger adversarial perturbations.

Applying Theorem 5.1 with flipping distribution $\mathcal{F}(\mathbf{x}, \boldsymbol{\theta})$ to each of the D softmax vectors of our model’s outputs yields l_0 -norm certificates for binary data that can be computed in linear time (see § F.3.1). In § F.3.2, we also apply it to the sparsity-aware smoothing distribution (Bojchevski et al., 2020), allowing us to differentiate between adversarial deletions and additions of bits. Theorem 5.1 can also be generalized to continuous distributions (see § F.3.3). But, for fair comparison with our baselines, we use the certificates of Eiras et al. (2022) as our base certificates for continuous data. In practice, the smoothed classifier and the base certificates cannot be evaluated exactly. One has to use Monte Carlo sampling to provide guarantees that hold with high probability (see § G).

6 LIMITATIONS

A limitation of our approach is that it assumes soft locality. It can be applied to arbitrary models, but may not necessarily result in better certificates than isotropic smoothing (recall § 4.3). Also, choosing the smoothing distributions requires some assumptions about which parts of the input are how relevant to making a prediction. Our experiments show that natural assumptions like homophily can be sufficient. But choosing a distribution may be more challenging for other tasks. A limitation of (most) randomized smoothing methods is that they use sampling to approximate the smoothed classifier. Because we use multiple distributions, we can only use a fraction of the samples per distribution. We can alleviate this problem by sharing smoothing distributions among outputs (see § E.1). Still, future work should try to improve the sample efficiency of randomized smoothing or develop deterministic base certificates (e.g. by generalizing (Levine & Feizi, 2020) to anisotropic distributions), which could then be incorporated into our linear programming framework.

³This term is equivalent to the exponential of the Rényi-divergence $\exp(\mathcal{D}_{\alpha}(\Psi_{\mathbf{x}'} || \Psi_{\mathbf{x}}))$ with $\alpha = 2$.

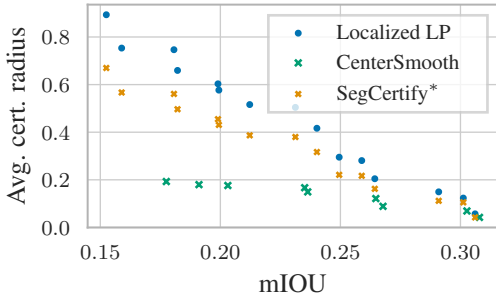


Figure 2: Comparison of isotropic smoothing with $\sigma_{\text{iso}} \in \{0.01, \dots, 0.5\}$ to our LP-based certificate with $(\sigma_{\text{min}}, \sigma_{\text{max}}) = (\sigma_{\text{iso}}, \infty)$, using a modified, strictly local U-Net on Pascal-VOC. Localized smoothing offers the same mIOU as SegCertify* and stronger robustness certificates.

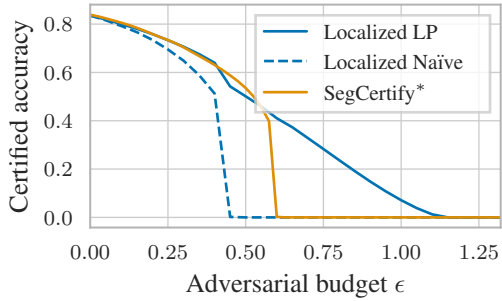


Figure 3: Certified accuracy of U-Net on Pascal-VOC. We compare SegCertify* ($\sigma_{\text{iso}} = 0.2$) to localized smoothing ($(\sigma_{\text{min}}, \sigma_{\text{max}}) = (0.15, 1.0)$). Combining the base certificates (dashed blue line) via our collective LP (solid blue line) outperforms the baseline.

7 EXPERIMENTAL EVALUATION

In this section, we compare our method to all existing collective certificates for ℓ_p -norm perturbations: Center smoothing using isotropic Gaussian noise (Kumar & Goldstein, 2021), SegCertify (Fischer et al., 2021) and the collective certificates of Schuchardt et al. (2021). To compare SegCertify to the other methods, we report the number of certifiably robust predictions and not just whether all predictions are robust. We write SegCertify* to highlight this. When considering models that are not strictly local (i.e. all outputs depend on all inputs) the certificates of Schuchardt et al. (2021) and Fischer et al. (2021) are identical, i.e., do not have to be evaluated separately. A more detailed description of the experimental setup, hardware and computational cost can be found in § C.

Metrics. Evaluating randomized smoothing methods based on certificate strength alone is not sufficient. Different distributions lead to different tradeoffs between prediction quality and certifiable robustness (as discussed in § 4.3). As metrics for prediction quality, we use *accuracy* and *mean intersection over union* (mIOU).⁴ The main metric for certificate strength is the *certified accuracy* $\xi(\epsilon)$, i.e., the percentage of predictions that are correct and certifiably robust, given adversarial budget ϵ . Following (Schuchardt et al., 2021), we use the *average certifiable radius* (ACR) as an aggregate metric, i.e. $\sum_{n=1}^{N-1} \epsilon_n \cdot (\xi(\epsilon_n) - \xi(\epsilon_{n+1}))$ with budgets $\epsilon_1 \leq \epsilon_2 \leq \dots \leq \epsilon_N$ and $\epsilon_1 = 0, \xi(\epsilon_N) = 0$.

Evaluation procedure. We assess the accuracy-robustness tradeoff of each method by computing accuracy / mIOU and ACR for a wide range of smoothing distribution parameters. We then eliminate all points that are Pareto-dominated, i.e. for which there exist different parameter values that yield higher accuracy / mIOU and ACR. Finally, we assess to if localized smoothing dominates the baselines, i.e. whether it can be parameterized to achieve strictly better accuracy-robustness tradeoffs.

7.1 IMAGE SEGMENTATION

Dataset and model. We evaluate our certificate for ℓ_2 perturbations on 100 images from the Pascal-VOC (Everingham et al., 2010) 2012 segmentation validation set. Training is performed on 10582 samples extracted from SBD, also known as "Pascal trainaug" (Hariharan et al., 2011). Additional experiments on Cityscapes (Cordts et al., 2016) can be found in § A. To increase batch sizes and thus allow a thorough investigation of different smoothing parameters, all images are downsampled to 50% of their original size, similar to (Fischer et al., 2021). Our base model is a U-Net segmentation model (Ronneberger et al., 2015) with a ResNet-18 backbone. For isotropic randomized smoothing, we use Gaussian noise $\mathcal{N}(0, \sigma_{\text{iso}})$ with different $\sigma_{\text{iso}} \in \{0.01, 0.02, \dots, 0.5\}$. To perform localized randomized smoothing, we choose parameters $\sigma_{\text{min}}, \sigma_{\text{max}} \in \mathbb{R}_+$ and partition all images into regular grids (similar to Fig. 1). To smooth outputs in grid cell (i, j) , we sample noise for grid cell (k, l) from $\mathcal{N}(0, \sigma' \cdot \mathbf{1})$, with $\sigma' \in [\sigma_{\text{min}}, \sigma_{\text{max}}]$ chosen proportional to the distance of (i, j) and

⁴I.e. add up confusion matrices over the entire dataset, compute per-class IOUs and average over all classes.

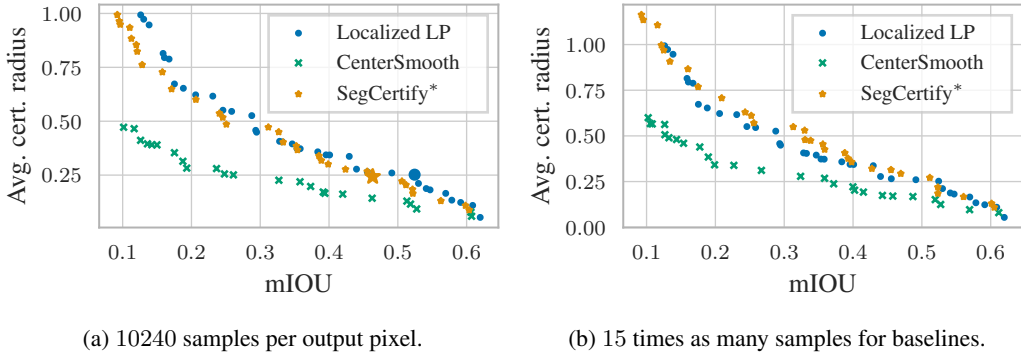


Figure 4: Comparison of isotropic smoothing to our LP-based certificate with a 3×5 grid and U-Net on Pascal-VOC. U-Net is sufficiently local to benefit from localized smoothing (Fig. 4a), but not enough to offset the increased sample complexity (Fig. 4b) for the probabilistic base certificates.

(k, l) (more details in § C.2). All training data is randomly perturbed using samples from the same smoothing distribution that is used for certification.

Accuracy-robustness tradeoff under strict locality. Our goal is to verify that, if a model is sufficiently local, localized smoothing offers a better accuracy-robustness tradeoff than isotropic smoothing. As an extreme example, we construct a strictly local model from our U-Net segmentation model. This modified model partitions each image into a grid of size 2×2 . It then iterates over cells (i, j) , sets all values outside (i, j) to 0 and applies the original model. Finally, it stitches all 4 segmentation masks into a single one. For such a strictly local model, we can apply localized smoothing with the same 2×2 grid and $\sigma_{\max} \rightarrow \infty$ to recover the certificate of Schuchardt et al. (2021) (see § I).⁵ Fig. 2 compares the resulting trade-off for $\sigma_{\min} = \sigma_{\text{iso}}$ to that of both isotropic smoothing baselines using 153600 Monte Carlo samples. Localized smoothing yields the same mIOUs as SegCertify*, but up to 22.4 p.p. larger ACR. Both approaches Pareto-dominate center smoothing.

Accuracy-robustness tradeoff under soft locality. Next, we want to verify our claim about the existence of softly local models for which localized smoothing is beneficial. To this end, we randomly smooth the U-Net model itself, without using masking to enforce strict locality. We perform localized smoothing with grid size 3×5 , various $\sigma_{\min} \in \{0.01, 0.02, \dots, 0.5\}$, $\sigma_{\max} \in [0.02, 1.0]$ and 10240 samples per output pixel (i.e. $10240 \cdot 15 = 153600$ samples in total). Isotropic smoothing is also performed with 10240 samples per output pixel. Fig. 4a shows that localized smoothing Pareto-dominates SegCertify* for high-accuracy models with $\text{mIOU} > 35.3\%$. Importantly the figure is not to be read like a line graph! Even if the vertical distance between two methods is small, one may significantly outperform the other. For example, $\sigma_{\text{iso}} = 0.1$, with an mIOU of 46.34% and an ACR of 0.24 (highlighted with a bold cross) is dominated by $(\sigma_{\min}, \sigma_{\max}) = (0.09, 0.2)$ (highlighted with a large circle), which has a larger ACR of 0.25 and a mIOU that is a whole 6.1 p.p. higher.

Benefit of linear programming. Fig. 3 demonstrates how the linear program derived in § 4.1 enables this improved tradeoff. We compare SegCertify* with $\sigma_{\text{iso}} = 0.2$ to localized smoothing with $(\sigma_{\min}, \sigma_{\max}) = (0.15, 1.0)$. Naïvely combining the base certificates (dashed line) is not sufficient for outperforming the baseline, as they cannot certify robustness beyond $\epsilon = 0.45$. However, solving the collective LP (solid blue line) extends the maximum certifiable radius to $\epsilon = 1.15$.

Sample efficiency. Using the same number of samples per output pixel for both localized and isotropic smoothing neglects that localized smoothing requires sampling from 15 different distributions, i.e. sampling 15 times as many images.⁶ In Fig. 4b we allow the baselines to sample the same number of images. Now, localized smoothing is mostly dominated by SegCertify*, except for high-accuracy models with $\text{mIOU} \in [52.4\%, 57.8\%]$ or $\text{mIOU} > 60.8\%$. We conclude that U-Net is local enough to benefit from localized smoothing, but not enough to offset the practical problem of having to work with fewer Monte Carlo samples (see also discussion in § 6) in the entire range of possible isotropic smoothing parameters. Note, however, that we can always recover the guarantees of SegCertify* by using a 1×1 grid (see § H).

⁵Note that they never evaluated their approach on image segmentation.

⁶This is however not necessary for the previously discussed strictly local model.

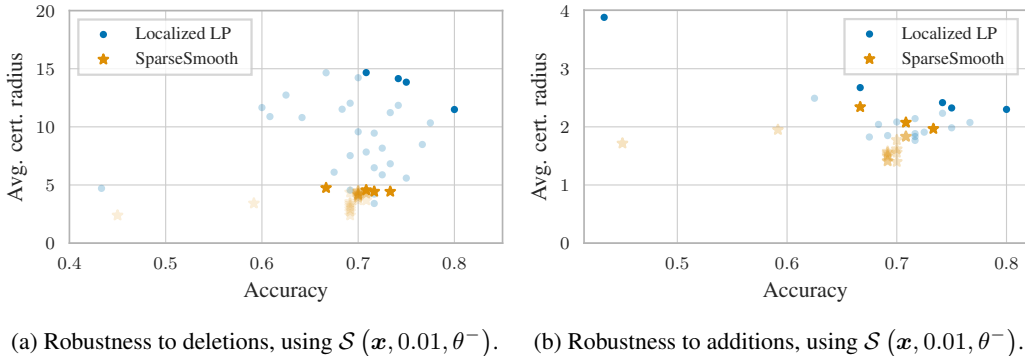


Figure 5: Comparison of our LP-based collective certificate to Bojchevski et al. (2020), using APPNP on Citeseer. We consider both adversarial deletions (Fig. 5a) and additions (Fig. 5b) of attribute bits. Locally smoothed models offer a better accuracy-robustness tradeoff, especially for deletions. Transparent points signal that they are Pareto-dominated by points from the same method.

7.2 NODE CLASSIFICATION ON CITSEER

Dataset and model. Finally, we consider models that are designed with locality in mind: Graph neural networks. We take APPNP (Klicpera et al., 2019), which aggregates per-node predictions from the entire graph based on personalized pagerank scores, and apply it to the Citeseer (Sen et al., 2008) dataset. To certify its robustness, we perform randomized smoothing with sparsity-aware noise $\mathcal{S}(\mathbf{x}, \theta^+, \theta^-)$, where θ^+ and θ^- control the probability of randomly adding or deleting node attributes, respectively (more details in § F.3.2). As a baseline we apply the tight certificate SparseSmooth of Bojchevski et al. (2020) to distributions $\mathcal{S}(\mathbf{x}, 0.01, \theta_{\text{iso}}^-)$ with $\theta_{\text{iso}}^- \in \{0.1, 0.15, \dots, 0.95\}$. The small addition probability 0.01 is meant to preserve the sparsity of the graph’s attribute matrix and was used in most experiments in (Bojchevski et al., 2020). For localized smoothing, we partition the graph into 5 clusters and define a minimum deletion probability $\theta_{\text{min}}^- \in \{0.1, 0.15, \dots, 0.95\}$. We then sample each cluster’s attributes from $\mathcal{S}(\mathbf{x}, 0.01, \theta'^-)$ with $\theta'^- \in [\theta_{\text{min}}^-, 0.95]$ chosen based on cluster affinity. To compute the base certificates, we use the variance-constrained certificate from § F.3.2. In all cases, we take $5 \cdot 10^5$ samples (i.e. 10^5 per cluster for localized smoothing). Further discussions, as well as experiments on different models and datasets can be found in § B.

Accuracy-robustness tradeoff. Fig. 5 shows the accuracy and ACR pairs achieved by the naïve isotropic smoothing certificate and the LP-based certificate for localized smoothing. Despite having fewer samples per prediction, our method outperforms the baseline, offering higher accuracy certifying larger ACRs, especially for attribute deletions. Notably, in some cases, our approach even improves accuracy by over 7 p.p. percentage points compared to isotropically smoothed models. Similar to the observation made by Bojchevski et al. (2020) in their Section K, we also find that increasing the probability of attribute perturbations can improve accuracy to some extent. We posit that localized smoothing can leverage this phenomenon as a form of test-time regularization while preserving the crucial attributes of nearby nodes. In § B.1 we show that the stems from the smoothing scheme and is not solely due to using our novel variance-constrained certificate.

8 CONCLUSION

We proposed a novel approach to achieve provable collective robustness in multi-output classifiers that extends beyond strict locality, utilizing our introduced localized randomized smoothing scheme. Our approach involves smoothing different outputs with anisotropic smoothing distributions that match the model’s soft locality. We demonstrated how per-output certificates obtained through localized smoothing can be combined into a strong collective robustness certificate using (mixed-integer) linear programming. Our experiments indicate that localized smoothing can achieve superior accuracy-robustness tradeoffs compared to isotropic smoothing methods. However, not all models match our distance-based locality assumption, particularly for image segmentation tasks. Node classification tasks are more amenable to localized smoothing due to their inherent locality. Our results highlight the importance of locality in achieving collective robustness and emphasize the need for future research to develop effective local models for multi-output tasks.

9 REPRODUCIBILITY STATEMENT

We prove all theoretic results that were not already derived in the main text in § D to § G. To ensure reproducibility of the experimental results we provide detailed descriptions of the evaluation process with the respective parameters in § C. An implementation, including configuration files, will be made available at <https://www.cs.cit.tum.de/daml/localized-smoothing>.

10 ETHICS STATEMENT

In this paper, we propose a method to increase the robustness of machine learning models against adversarial perturbations and to certify their robustness. We see this as an important step towards general usage of models in practice, as many existing methods are brittle to crafted attacks. Through the proposed method, we hope to contribute to the safe usage of machine learning. However, robust models also have to be seen with caution. As they are harder to fool, harmful purposes like mass surveillance are harder to avoid. We believe that it is still necessary to further research robustness of machine learning models as the positive effects can outweigh the negatives, but it is necessary to discuss the ethical implications of the usage in any specific application area.

11 ACKNOWLEDGEMENTS

This research is funded by the Bavarian Ministry of Economic Affairs, Regional Development and Energy with funds from the Hightech Agenda Bayern. Further, it is supported by the German Research Foundation, grant GU 1409/4-1.

REFERENCES

- T.W. Anderson. Confidence limits for the value of an arbitrary bounded random variable with a continuous distribution function. In *Bulletin of The International and Statistical Institute*, volume 43, pp. 249–251, 1969.
- Yonatan Belinkov and Yonatan Bisk. Synthetic and natural noise both break neural machine translation. In *International Conference on Learning Representations*, 2018.
- Aleksandar Bojchevski and Stephan Günnemann. Deep gaussian embedding of graphs: Unsupervised inductive learning via ranking. In *International Conference on Learning Representations*, 2018.
- Aleksandar Bojchevski and Stephan Günnemann. Certifiable robustness to graph perturbations. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. Efficient robustness certificates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more. In *International Conference on Machine Learning*, 2020.
- Gregory Bonaert, Dimitar I. Dimitrov, Maximilian Baader, and Martin Vechev. Fast and precise certification of transformers. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, 2021.
- Carlo Bonferroni. Teoria statistica delle classi e calcolo delle probabilita. *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze*, 8:3–62, 1936.
- G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000.
- Alexander Buslaev, Vladimir I. Iglovikov, Eugene Khvedchenya, Alex Parinov, Mikhail Druzhinin, and Alexandr A. Kalinin. Alumentations: Fast and flexible image augmentations. *Information*, 11, 2020.
- Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation, 2017.

- Ruoxin Chen, Zenan Li, Jie Li, Junchi Yan, and Chentao Wu. On collective robustness of bagging against data poisoning. In *Proceedings of the 39th International Conference on Machine Learning*, 2022.
- C. J. Clopper and E. S. Pearson. The use of confidence or fiducial limits illustrated in the case of the binomial. *Biometrika*, 26:404–413, 1934.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 2016.
- A. Dvoretzky, J. Kiefer, and J. Wolfowitz. Asymptotic Minimax Character of the Sample Distribution Function and of the Classical Multinomial Estimator. *The Annals of Mathematical Statistics*, 27:642 – 669, 1956.
- Francisco Eiras, Motasem Alfarra, Philip Torr, M. Pawan Kumar, Puneet K. Dokania, Bernard Ghanem, and Adel Bibi. ANCER: Anisotropic certification via sample-wise volume maximization. *Transactions on Machine Learning Research*, 2022.
- Mark Everingham, Luc Gool, Christopher K. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vision*, 88:303–338, June 2010.
- Marc Fischer, Maximilian Baader, and Martin Vechev. Certified defense to image transformations via randomized smoothing. In *Advances in Neural Information Processing Systems*, 2020.
- Marc Fischer, Maximilian Baader, and Martin Vechev. Scalable certified segmentation via randomized smoothing. In *International Conference on Machine Learning*, 2021.
- Manuel Gil, Fady Alajaji, and Tamas Linder. Rényi divergence measures for commonly used univariate continuous distributions. *Information Sciences*, 249:124–131, 2013.
- Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *2011 International Conference on Computer Vision*, 2011.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Sture Holm. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6(2):65–70, 1979.
- George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1998.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations*, 2017.
- Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate: Graph neural networks meet personalized pagerank. In *International Conference on Learning Representations*, 2019.
- Ching-Yun Ko, Zhaoyang Lyu, Lily Weng, Luca Daniel, Ngai Wong, and Dahua Lin. POPQORN: Quantifying robustness of recurrent neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Aounon Kumar and Tom Goldstein. Center smoothing: Provable robustness for functions with metric-space outputs. In *Advances in Neural Information Processing Systems*, 2021.

- Aounon Kumar, Alexander Levine, Soheil Feizi, and Tom Goldstein. Certifying confidence via randomized smoothing. In *Advances in Neural Information Processing Systems*, 2020.
- Mathias Lécuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy*, pp. 656–672, 2019.
- Guang-He Lee, Yang Yuan, Shiyu Chang, and Tommi S. Jaakkola. Tight certificates of adversarial robustness for randomly smoothed classifiers. In *Advances in Neural Information Processing Systems*, 2019.
- Alexander Levine and Soheil Feizi. (de)randomized smoothing for certifiable defense against patch attacks. In *Advances in Neural Information Processing Systems*, 2020.
- Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive noise. In *Advances in Neural Information Processing Systems*, 2019.
- Yongge Liu, Jianzhuang Yu, and Yahong Han. Understanding the effective receptive field in semantic image segmentation. *Multimedia Tools and Applications*, 77(17):22159–22171, 2018.
- Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016.
- Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the construction of internet portals with machine learning. *Information Retrieval*, 2000.
- MOSEK ApS. *MOSEK Optimizer API for Python 9.2.46*, 2019. URL <https://docs.mosek.com/9.2/pythonapi/index.html>.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, 2015.
- Wonryong Ryou, Jiayu Chen, Mislav Balunovic, Gagandeep Singh, Andrei Dan, and Martin Vechev. Scalable polyhedral verification of recurrent neural networks. In Alexandra Silva and K. Rustan M. Leino (eds.), *Computer Aided Verification*, 2021.
- Jan Schuchardt, Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. Collective robustness certificates: Exploiting interdependence in graph neural networks. In *International Conference on Learning Representations*, 2021.
- Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI Magazine*, 29:93, 2008.
- Han Shi, Jiahui Gao, Xiaozhe Ren, Hang Xu, Xiaodan Liang, Zhenguo Li, and James Tin-Yau Kwok. Sparsebert: Rethinking the importance analysis in self-attention. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, 2021.
- Zhouxing Shi, Huan Zhang, Kai-Wei Chang, Minlie Huang, and Cho-Jui Hsieh. Robustness verification for transformers. In *International Conference on Learning Representations*, 2020.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- Hoang-Dung Tran, Neelanjana Pal, Patrick Musau, Diego Manzananas Lopez, Nathaniel Hamilton, Xiaodong Yang, Stanley Bak, and Taylor T. Johnson. Robustness verification of semantic segmentation neural networks using relaxed reachability. In *Computer Aided Verification*, 2021.
- Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. Adversarial examples for semantic segmentation and object detection. In *International Conference on Computer Vision*, 2017.

- Pavel Yakubovskiy. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2020.
- Maksym Yatsura, Kaspar Sakmann, N. Grace Hua, Matthias Hein, and Jan Hendrik Metzen. Certified defences against adversarial patch attacks on semantic segmentation. In *Workshop on Trustworthy and Socially Responsible Machine Learning, NeurIPS 2022*, 2022.
- Dinghui Zhang, Mao Ye, Chengyue Gong, Zhanxing Zhu, and Qiang Liu. Black-box certification with randomized smoothing: A functional optimization based framework. In *Advances in Neural Information Processing Systems*, 2020.
- Daniel Zügner and Stephan Günnemann. Adversarial attacks on graph neural networks via meta learning. In *International Conference on Learning Representations*, 2019.
- Daniel Zügner and Stephan Günnemann. Certifiable robustness and robust training for graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019.
- Daniel Zügner and Stephan Günnemann. Certifiable robustness of graph convolutional networks under structure perturbations. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020.

A	Image Segmentation on CityScapes	15
B	Additional Experiments on Node Classification	16
B.1	Comparison to the Naïve Variance-constrained Isotropic Smoothing Certificate . . .	16
B.2	Node Classification using Graph Convolutional Networks	16
B.3	Node classification on Cora-ML	17
B.4	Lower Certifiable Robustness to Additions	18
B.5	Benefit of Linear Programming Certificates	19
C	Detailed Experimental Setup	21
C.1	Certificate Strength Metrics	21
C.2	Semantic Segmentation	21
C.3	Node Classification	22
C.4	Hardware and runtime	23
D	Proof of Theorem 4.2	24
E	Improving Efficiency	25
E.1	Sharing Smoothing Distributions Among Outputs	25
E.2	Quantizing Certificate Parameters	25
E.3	Sharing Noise Levels Among Inputs	26
E.4	Linear Relaxation	27
F	Base Certificates	28
F.1	Gaussian Smoothing for l_2 Perturbations of Continuous Data	28
F.2	Uniform Smoothing for l_1 Perturbations of Continuous Data	29
F.3	Variance-Constrained Certification	29
G	Monte Carlo Randomized Smoothing	38
G.1	Monte Carlo Base Certificates for Continuous Data	38
G.2	Monte Carlo Variance-Constrained Certification	39
G.3	Monte Carlo Center Smoothing	40
G.4	Multiple Comparisons Problem	41
H	Comparison to the Collective Certificate of Fischer et al. (2021)	42
I	Comparison to the Collective Certificate of Schuchardt et al. (2021)	43
I.1	The Collective Certificate	43
I.2	Proof of Subsumption	44

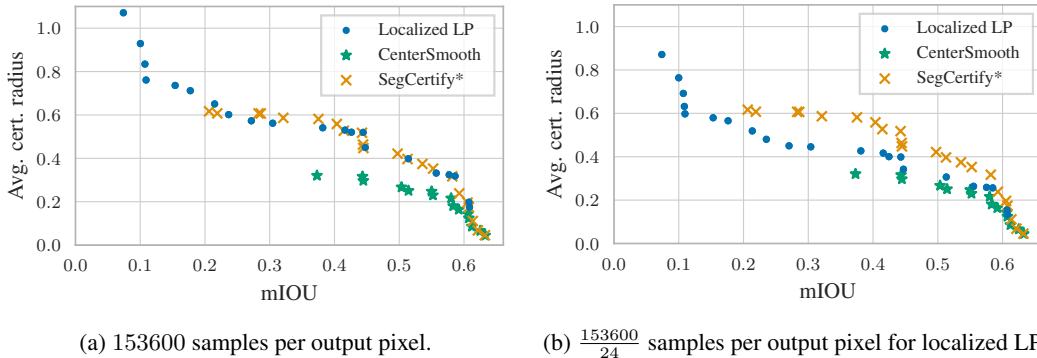


Figure 6: Comparison of our LP-based collective certificate for localized randomized smoothing with a 3×5 grid to CenterSmooth and SegCertify*, using DeepLabV3 on Cityscapes. Increasing the number of samples used for certifying each output from $6400 = \frac{153600}{24}$ to 153600 (same as for the baselines) closes the gap between localized randomized smoothing and SegCertify*. Still, localized smoothing only offers stronger certificates for models with $\text{mIOU} \leq 0.21$ (compared to $\text{mIOU} \leq 0.11$ when using fewer samples).

A IMAGE SEGMENTATION ON CITYSCAPES

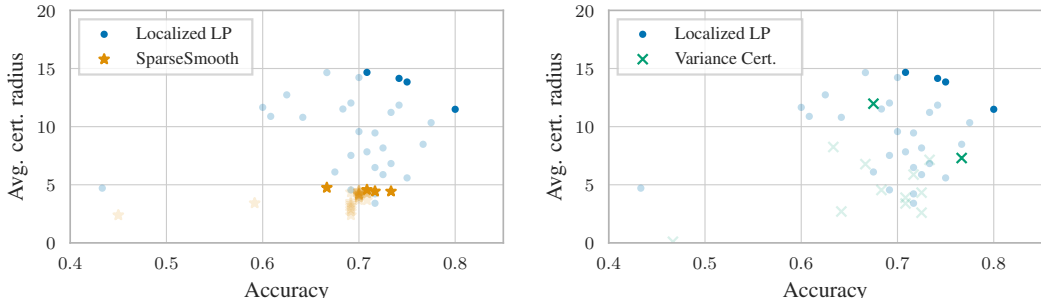
In the following, we apply our approach to DeepLabv3 (Chen et al., 2017) models trained on the Cityscapes (Cordts et al., 2016) training set. We evaluate the certificates on 50 images from the validation set. For localized smoothing, we partition the image into a grid of shape 4×6 . To limit the number of LP variables despite the increased resolution, we quantize the base certificate parameters $\eta^{(n)}$ into 2048 bins (see § E.2). Different from our experiments on Pascal-VOC and due to the increased computational cost of using higher-dimensional images, the locally smoothed models are not trained on the localized smoothing distribution with parameters $(\sigma_{\min}, \sigma_{\max})$. Instead, we use model trained with isotropic Gaussian noise with standard deviation $\sigma_{\text{iso}} = \sigma_{\min}$.

Fig. 6a shows that, even when allowing 153600 samples per output pixel for both localized smoothing and the baselines (i.e. localized smoothing gets to sample 24 times as many images), most choices of $(\sigma_{\min}, \sigma_{\max})$ do not offer higher accuracy and robustness than SegCertify*, except those leading to a small mIOU below 0.21. Fig. 6b shows that reducing the number of samples per output pixel for localized smoothing to $6400 = \frac{153600}{24}$ further weakens the certificate. There, localized smoothing only offers stronger certificates for models with an mIOU below 0.11.

There are three possible explanations for why localized smoothing does not outperform SegCertify*. The first one is that we do not train on the same distribution that we use for certification, so our models are less accurate or less consistent in their predictions, which reduces mIOU or certified robustness. The second one is that our simplistic choice of localized smoothing based on grid cell distance (see § C.2) does not match the actual locality structure of DeepLabv3. The last one is that DeepLabv3, which uses dilated convolutions to increase the receptive field size in each layer, is just inherently less local than the U-Net architecture used in our experiments on Pascal-VOC. Nevertheless, it should be noted that we can always parameterize localized smoothing to obtain the same results as SegCertify* (see Appendix H).

B ADDITIONAL EXPERIMENTS ON NODE CLASSIFICATION

In the following, we perform additional experiments on graph neural networks for node classification, including a different model and an additional dataset. Unless otherwise stated, all details of the experimental setup are identical to § 7.2. In particular, we use sparsity-aware smoothing distribution $\mathcal{S}(\mathbf{x}, 0.01, \theta^-)$, where probability of deleting bits θ^- is either constant across the entire graph (for the isotropic randomized smoothing baseline) or adjusted per output and cluster based on cluster affinity (for localized randomized smoothing).



(a) Using (Bojchevski et al., 2020) for the naïve isotropic smoothing baseline. (b) Using variance-constrained certification for the naïve isotropic smoothing baseline.

Figure 7: Analysis of our LP-based collective certificate using APPNP on Citeseer. We use the sparsity-aware smoothing with $\mathcal{S}(\mathbf{x}, 0.01, \theta^-)$ to certify robustness to deletions. In Fig. 7a we use the certificate of Bojchevski et al. (2020) for baseline (identical to Fig. 5a). In Fig. 7b we use variance-constrained certification (see Theorem 5.1) as baseline. In both cases, there are locally smoothed models with a higher accuracy than any of the isotropically smoothed models and significantly larger average certifiable radii.

B.1 COMPARISON TO THE NAÏVE VARIANCE-CONSTRAINED ISOTROPIC SMOOTHING CERTIFICATE

In Fig. 5 of § 7.2, we observed that locally smoothed models surprisingly did not only achieve up to three times higher average certifiable radii, but simultaneously had higher accuracy than any of the isotropically smoothed models. One potential explanation is that we used variance-constrained certification (see Theorem 5.1) (i.e. smoothing the models’ softmax scores instead of their predicted labels) for localized smoothing, but not for the isotropic smoothing baseline. This might result in two substantially different models. To investigate this, we repeat the experiment from Fig. 5a, using variance-constrained certification for both localized smoothing and the isotropic smoothing baseline. Fig. 7 shows that, no matter which smoothing paradigm we use for our isotropic smoothing baseline, there is a c.a. 7 p.p. difference in accuracy between the most accurate isotropically smoothed model and the most accurate locally smoothed model.

Interestingly, even variance-constrained smoothing with isotropic noise (green crosses in Fig. 7b) is sufficient for outperforming the isotropic smoothing certificate of Bojchevski et al. (2020) (orange stars in Fig. 7a). This showcases that variance-constrained certification does not only present a very efficient, but also a very effective way of certifying robustness on discrete data (even when entirely ignoring the collective robustness aspect).

B.2 NODE CLASSIFICATION USING GRAPH CONVOLUTIONAL NETWORKS

So far, we have only used APPNP models as our base classifier. Now, we repeat our experiments using 6-layer Graph Convolutional Networks (GCN) (Kipf & Welling, 2017). In each layer, GCNs first apply a linear layer to each node’s latent vector and then average over each node’s 1-hop neighborhood. Thus, a 6-layer GCN classifies each node using attributes from all nodes in its 6-hop neighborhood, which covers most or all of the Citeseer graph. Aside from using GCN instead of APPNP as the base model, we leave the experimental setup from § 7.2 unchanged. Note that GCNs

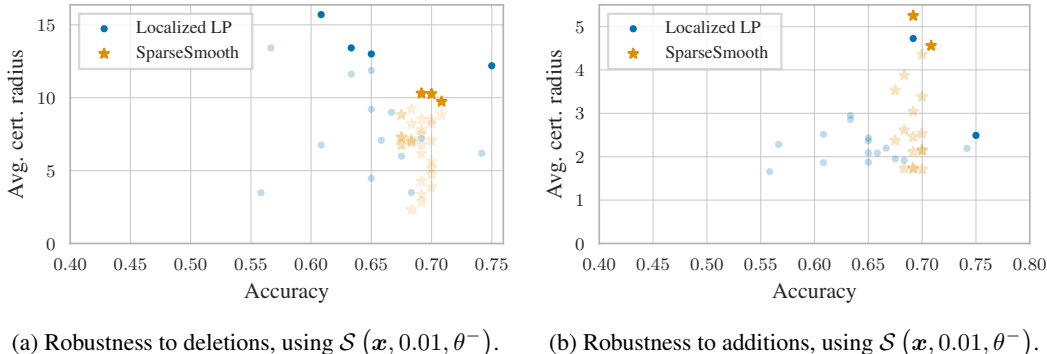


Figure 8: Comparison of our LP-based collective certificate for localized randomized smoothing to SparseSmooth, using a 6-layer GCN on Citeseer. We consider both adversarial deletions (Fig. 8a) and additions (Fig. 8b). Some locally smoothed models have a higher accuracy than any of the isotropically smoothed models. However, our certificate only dominates the best isotropically smoothed models when considering robustness to deletions, not when considering robustness to additions. This can either be attributed to a lower locality in deep GCNs or variance-constrained certification yielding weak base certificates for addition when θ^+ is small.

are typically used with fewer layers. However, these shallow models are strictly local and it has already been established that the certificate Schuchardt et al. (2021) – which is subsumed by our certificate (see § I.2) – can provide very strong robustness guarantees for them. We therefore increase the number of layers to obtain a model that is not strictly local.

Fig. 8 shows the results for both robustness to deletions and robustness to additions. Similar to APPNP, some locally smoothed models have an up to 4 p.p. higher accuracy than the most accurate isotropically smoothed model. When considering robustness to deletions, the locally smoothed models Pareto-dominate all of the isotropically smoothed models, i.e. offer better accuracy-robustness tradeoffs. Some can guarantee average certifiable radii that are at least 50% larger than those of the baseline. When considering robustness to additions however, some of the isotropically smoothed models have a higher certifiably robustness.

We see two potential causes for our method’s lower certifiable robustness to additions: The first potential cause is that the GCN may be less local than APPNP or that it has a different form of locality that does not match our clustering-based localized smoothing distributions. This appears plausible, as GCN averages uniformly over each neighborhood, whereas APPNP aggregates predictions based on pagerank scores. APPNP may thus primarily attend to specific, densely connected nodes, making it more local than GCN. The second potential cause is that the variance-constrained certificate we use as our base certificate may be less effective when certifying robustness to adversarial additions by using a very small addition probability like $\theta^+ = 0.01$. Afterall, we have also seen in our experiments with APPNP in § 7.2 that the gap in average certifiable radii between localized and isotropic smoothing was significantly smaller when considering additions. We investigate this second potential cause in more detail in § B.4.

B.3 NODE CLASSIFICATION ON CORA-ML

Next, we repeat our experiments with APPNP on the Cora-ML (McCallum et al., 2000; Bojchevski & Günnemann, 2018) node classification dataset, keeping all other parameters fixed. The results are shown in Fig. 9. Unlike on Citeseer, the locally smoothed models have a slightly reduced accuracy compared to the isotropically smoothed models. This can either be attributed to one smoothing approach having a more desirable regularizing effect on the neural network, or the fact that we smooth softmax scores instead of predicted labels when constructing the locally smoothed models. Nevertheless, when considering adversarial deletions, localized smoothing makes it possible to achieve average certifiable radii that are at least 50% larger than any of the isotropically smoothed models’ – at the cost of slightly reduced accuracy 8.6%. Or, for another point of the pareto front, we in-

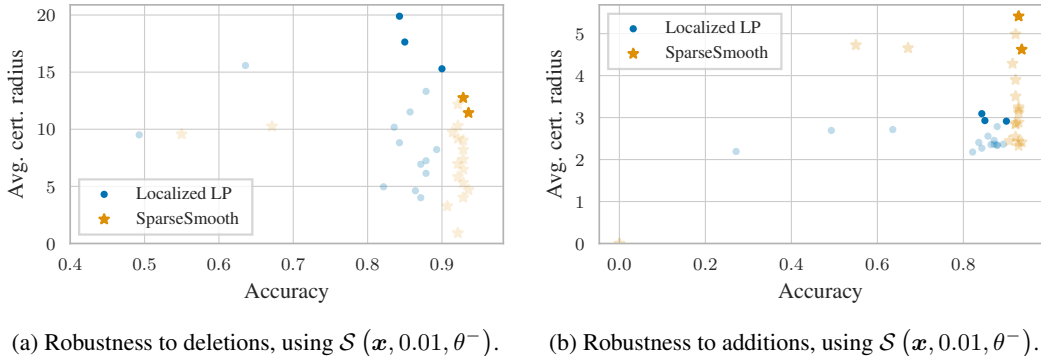


Figure 9: Comparison of our LP-based collective certificate for localized randomized smoothing to the SparseSmooth (Bojchevski et al., 2020), using APPNP on Cora-ML. We consider both adversarial deletions (Fig. 9a) and additions (Fig. 9b). Some locally smoothed models that have a higher accuracy than any of the isotropically smoothed models. However, our method is only able to dominate all isotropically smoothed models when considering robustness to deletions, not when considering robustness to additions. This can either be attributed to a lower locality in deep GCNs or variance-constrained certification yielding weak base certificates for addition when θ^+ is small.

crease the certificate by 20% while reducing the accuracy by 2.8 percentage points. As before, the certificates for attribute additions are significantly weaker.

B.4 LOWER CERTIFIABLE ROBUSTNESS TO ADDITIONS

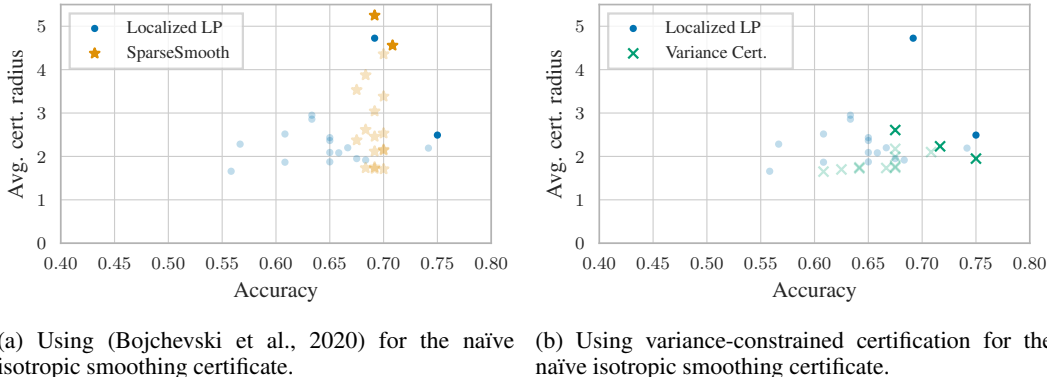


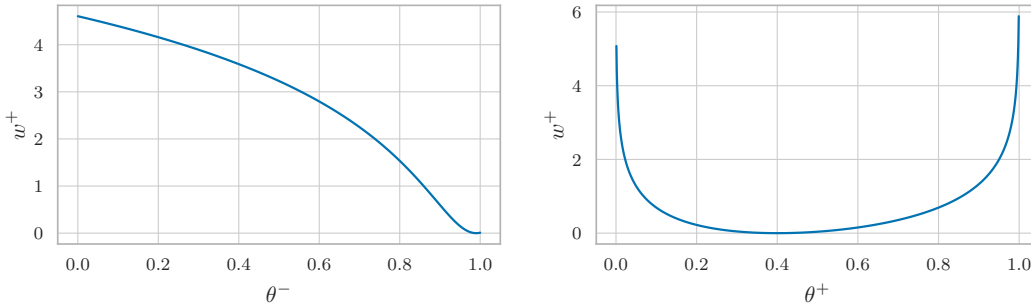
Figure 10: Comparison of our LP-based collective certificate for localized randomized smoothing to SparseSmooth and to a naïve combination of its base certificates, using GCN and adversarial additions on Citeseer. Fig. 10a shows that the LP-based certificate is outperformed by naïve isotropic smoothing. Fig. 10b shows that this is largely due to the variance-constrained base certificates (green crosses) for adversarial additions being much weaker than the isotropic smoothing certificate of (Bojchevski et al., 2020) in Fig. 10a.

While our certificates for adversarial deletions have compared favorably to the isotropic smoothing baseline in all previous experiments, our certificates for adversarial additions were comparatively weaker on Cora-ML and when using GCNs as base models. In the following, we investigate to what extent this can be attributed to our use of variance-constrained certification for our base certificates.

Fig. 10a shows both our linear programming collective certificate and the naïve isotropic smoothing certificate based on (Bojchevski et al., 2020) for GCNs on Citeseer under adversarial additions. In Fig. 10b, we plot not only the LP-based certificates, but also our variance-constrained base certificates (drawn as green crosses). Comparing both figures shows that our base certificate’s average certifiable radii are at least 50% smaller than the largest ACR achieved by (Bojchevski et al., 2020)

in Fig. 10a. While our linear program significantly improves upon them, it is not sufficient to overcome this significant gap. This result is in stark contrast to our results for attribute deletions § B.1, where the variance-constrained base certificates alone were enough to significantly outperform the certificate of (Bojchevski et al., 2020).

Now that we have established that the variance-constrained base certificates appear significantly weaker for additions, we can analyze why. For this, recall that our base certificates are parameterized by a weight vector w (see Definition 4.1), with smaller values corresponding to higher robustness – or two weight vectors w^+ , w^- quantifying robustness to adversarial additions and deletions, respectively (see § F.3.2). Using our results from § F.3.2, we can draw the weights w^+ resulting from smoothing distribution $\mathcal{S}(x, 0.01, \theta^-)$ as a function of θ^- . Fig. 11a shows that θ^- has to be brought very close to 1 in order to guarantee high robustness to deletions, effectively deleting almost all attributes in the graph. Alternatively, one can also increase the addition probability θ^+ to perhaps 10% or 20%. But this would utterly destroy the sparsity of the graph’s attribute matrix. We can conclude that, while variance-constrained certification can in principle provide strong certificates for attribute deletions, it might be a worse choice than the method of Bojchevski et al. (2020) for very sparse datasets that force the use of very low addition probabilities θ^+ .



(a) Certificate weight w^+ for $\mathcal{S}(x, 0.01, \theta^-)$ for varying θ^- . (b) Certificate weight w^+ for $\mathcal{S}(x, \theta^+, 0.6)$ for varying θ^+ .

Figure 11: Base certificate weight w^+ of the variance-constrained sparsity-aware smoothing certificate for varying distribution parameters. Certifying high robustness to adversarial additions (i.e. obtaining small weights) requires either setting a high probability for random additions or an even higher probability for random deletions.

B.5 BENEFIT OF LINEAR PROGRAMMING CERTIFICATES

As we did for our experiments on image segmentation (see Fig. 3), we can inspect the certified accuracy curves of specific smoothed models in more detail to gain a better understanding of how the collective linear programming certificate enables larger average certifiable radii. We use the same experimental setup as in § 7.2, i.e. APPNP on Citeseer, and certify robustness to deletions. We compare the certifiably most robust isotropically smoothed model ($\theta_{\text{iso}}^- = 0.8$, ACR = 5.67 to the locally smoothed model with $\theta_{\text{min}}^- = 0.75$, $\theta_{\text{max}}^+ = 0.95$. For the locally smoothed models, we compute both LP-based collective certificate, as well as the naïve collective certificate.

Fig. 12 shows that even naïvely combining the localized smoothing base certificates obtained via variance-constrained certification (dashed blue line) is sufficient for outperforming the naïve isotropic smoothing certificate. This speaks to its effectiveness as a certificate against adversarial deletions. Combining the base certificates via linear programming (solid blue line) significantly enlarges this gap, leading to even larger maximum and average certifiable radii.

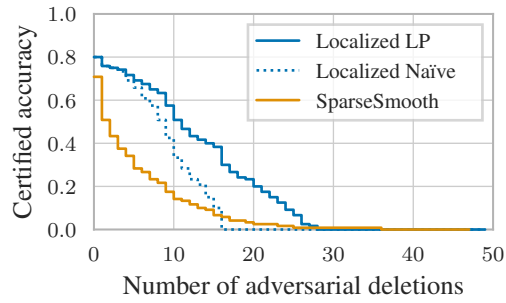


Figure 12: Certified accuracy of APPNP on Citeseer. We compare the naïve isotropic smoothing certificate of the most robust baseline model ($\theta_{\text{iso}}^- = 0.8$) to localized smoothing ($\theta_{\text{min}}^- = 0.75$). Even naïvely combining the variance-constrained base certificates (dashed blue line) is sufficient for outperforming the SparseSmooth certificate for 15 deletions or less. Combining the base certificates via our LP (solid blue line) further extends the certifiable radius and significantly increases the certified accuracy for perturbations with 5 or more deletions.

C DETAILED EXPERIMENTAL SETUP

In the following, we first explain the metrics we use for measuring the strength of certificates, and how they can be applied to the different types of randomized smoothing certificates used in our experiments. We then discuss the specific parameters and hyperparameters for our semantic segmentation and node classification experiments. We conclude by specifying the used hardware and comparing the computational cost of Monte Carlo sampling to that of solving the collective linear program.

C.1 CERTIFICATE STRENGTH METRICS

We use two metrics for measuring certificate strength: For specific adversarial budgets ϵ , we compute the certified accuracy $\xi(\epsilon)$ (i.e. the percentage of correct and certifiably robust predictions). As an aggregate metric, we compute the average certifiable radius, i.e. the lower Riemann integral of $\xi(\epsilon)$ evaluated at $\epsilon_1, \dots, \epsilon_N$ with $\epsilon_1 = 0$ and $\xi(\epsilon_N) = 0$. For our experiments on image segmentation, we use 81 equidistant points in $[0, 4]$. For our experiments on node classification, where we certify robustness to a discrete number of perturbations, we use $\epsilon_n = n$, i.e. natural numbers. In all experiments, we perform Monte Carlo randomized smoothing (see § G). Therefore, we may have to abstain from making predictions. Abstentions are counted as non-robust and incorrect. In the case of center smoothing, either all or no predictions abstain (this is inherent to the method. In our experiments, center smoothing never abstained).

C.1.1 COMPUTING CERTIFIED ACCURACY

The three different types of collective certificate considered in our experiments each require a different procedure for computing the certified accuracy. In the following, let $\mathbb{Z} = \{d \in \{1, \dots, D_{\text{out}}\} \mid f_n(\mathbf{x}) = \hat{y}_n\}$ be the indices of correct predictions, given an input \mathbf{x} .

Naïve collective certificate. The naïve collective certificate certifies each prediction independently. Let $\mathbb{H}^{(n)}$ be the set of perturbed inputs y_n is certifiably robust to (see Definition 3.1). Let \mathbb{B}_x be the collective perturbation model. Then $\mathbb{L} = \{d \in \{1, \dots, D_{\text{out}}\} \mid \mathbb{B}_x \subseteq \mathbb{H}^{(n)}\}$ is the set of all certifiably robust predictions. The certified accuracy can be computed as $\frac{|\mathbb{L} \cap \mathbb{Z}|}{D_{\text{out}}}$.

Center smoothing Center smoothing used for collective robustness certification does not determine which predictions are robust, but only the number of robust predictions. We therefore have to make the worst-case assumption that the correct predictions are the first to be changed by the adversary. Let l be the number of certifiably robust predictions. The certified accuracy can then be computed as $\frac{\max(0, |\mathbb{Z}| - (D_{\text{out}} - l))}{D_{\text{out}}}$.

Collective certificate. Let $l(\mathbb{T})$ be the optimal value of our collective certificate for the set of targeted nodes \mathbb{T} . Then the certified accuracy can be computed via $\frac{l(\mathbb{T})}{D_{\text{out}}}$ with $\mathbb{T} = \mathbb{Z}$.

C.2 SEMANTIC SEGMENTATION

Here, we provide all parameters of our experiments on image segmentation.

Models. As base models for the semantic segmentation tasks, we use U-Net (Ronneberger et al., 2015) and DeepLabv3 (Chen et al., 2017) segmentation heads with a ResNet-18 (He et al., 2016) backbone, as implemented by the Pytorch Segmentation Models library (version 0.13) (Yakubovskiy, 2020). We use the library’s default parameters. In particular, the inputs to the U-Net segmentation head are the features of the ResNet model after the first convolutional layer and after each ResNet block (i.e. after every fourth of the subsequent layers). The U-Net segmentation head uses (starting with the original resolution) 16, 32, 64, 128 and 256 convolutional filters for processing the features at the different scales. For the DeepLabv3 segmentation head, we use all default parameters from Chen et al. (2017) and an output stride of 16. To avoid dimension mismatches in the segmentation head, all input images are zero-padded to a height and width that is the next multiple of 32.

Data and preprocessing. We evaluate our certificates on the Pascal-VOC 2012 and Cityscapes segmentation validation set. We do not use the test set, because evaluating metrics like the certified

accuracy requires access to the ground-truth labels. For training the U-Net models on Pascal, we use the 10582 Pascal segmentation masks extracted from the SBD dataset (Hariharan et al., 2011) (referred to as "Pascal trainaug" or "Pascal augmented training set" in other papers). SBD uses a different data split than the official Pascal-VOC 2012 segmentation dataset. We avoid data leakage by removing all training images that appear in the validation set. For training the DeepLabv3 model on Cityscapes, we use the default training set. We downscale both the training and the validation images and ground-truth masks to 50% of their original height and width, so that we can use larger batch sizes and thus use our compute time to more thoroughly evaluate a larger range of different smoothing distributions. The segmentation masks are downscaled using nearest-neighbor interpolation, the images are downscaled using the INTER_AREA operation implemented in OpenCV (Bradski, 2000).

Training and data augmentation. We initialize our model weights using the weights provided by the Pytorch Segmentation Models library, which were obtained by pre-training on ImageNet. We train our models for 512 epochs, using Dice loss and Adam($lr = 0.001, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}, \text{weight_decay} = 0$). We use a batch size of 128 for Pascal-VOC and a batch size of 32 for Cityscapes. Every 8 epochs, we compute the mean IOU on the validation set. After training, we use the model that achieved the highest validation mean IOU. We apply the following train-time augmentations: With 50% probability, each image is randomly scaled by a factor from $[1, 2.0]$ using the ShiftScaleRotate augmentation implemented by the Albumentations library (version 0.5.2) (Buslaev et al., 2020). The images are then cropped to a fixed size of 160×256 (for Pascal-VOC) or 384×384 (for Cityscapes). Where necessary, the images are padded with zeros. Padded parts of the segmentation mask are ignored by the loss function. After these operations, each input is randomly perturbed using Gaussian noise. For isotropic smoothing, we use a fixed standard deviation $\sigma_{\text{iso}} \in \{0, 0.01, \dots, 0.5\}$, i.e. we train 51 different models on different isotropic smoothing distributions. For localized smoothing with grid shape $H \times W$ and parameters $(\sigma_{\text{min}}, \sigma_{\text{max}})$ we perform localized smoothing with a single sample per image. Since this generates $H \cdot W$ as many perturbed images, we perform gradient accumulation, processing $\frac{1}{H \cdot W}$ of each batch at a time. All samples are clipped to $[0, 1]$ to retain valid RGB-values.

Certification. For Pascal-VOC, we evaluate all certificates on the first 100 images from the validation set that – after downscaling – have a resolution of 166×250 . For Cityscapes, we use every tenth image from the validation set. For all certificates, we use Monte Carlo randomized smoothing (see discussion in § G). We use the significance parameter α to 0.01, i.e. all certificates hold with probability 0.99. For the center smoothing baseline, we use the default parameters suggested by the authors ($\Delta = 0.05, \beta = 2, \alpha_1 = \alpha_2$). For the naïve isotropic randomized smoothing baseline and for localized smoothing, we use Holm correction to account for the multiple comparisons problem, which yields strictly better results than Bonferroni correction (see § G.4). For our localized smoothing distribution, we partition the input image into a regular grid of size $H \times W$ (specified in the different paragraphs of § 7.1) and define minimum standard deviation σ_{min} and maximum standard deviation σ_{max} . Let $\mathbb{J}^{(k,l)}$ be the set of all pixel coordinates in grid cell (k, l) . To smooth outputs in grid cell (i, j) , we use a smoothing distribution $\mathcal{N}(\mathbf{0}, \text{diag}(\boldsymbol{\sigma}))$ with $\forall k \in \{1, \dots, H\}, l \in \{1, \dots, W\}, d \in \mathbb{J}^{(k,l)}$,

$$\sigma_d = \sigma_{\text{min}} + (\sigma_{\text{max}} - \sigma_{\text{min}}) \cdot \frac{\max(|i - k|, |l - j|)}{W}, \quad (7)$$

i.e. we linearly interpolate between σ_{min} and σ_{max} based on the l_∞ distance of grid cells (i, j) and (k, l) . All results are reported for the relaxed linear programming formulation of our collective certificate (see § E.4). The collective linear program is solved using MOSEK (version 9.2.46) (MOSEK ApS, 2019) through the CVXPY interface (version 1.1.13)

C.3 NODE CLASSIFICATION

Here, we provide all parameters of our experiments on node classification.

Model We test two different models: 2-layer APPNP (Klicpera et al., 2019) and 6-layer GCN (Kipf & Welling, 2017). For both models we use a hidden size of 64 and dropout with a probability of 0.5. For the propagation step of APPNP we use 10 for the number of iterations and 0.15 as the teleport probability.

Data and preprocessing. We evaluate our approach on the Cora-ML and Citeseer node classification datasets. We perform standard preprocessing, i.e., remove self-loops, make the graph undirected and select the largest connected component. We use the same data split as in (Schuchardt et al., 2021), i.e. 20 nodes per class for the train and validation set.

Training and data augmentation All models are trained with a learning rate of 0.001 and weight decay of 0.001. The models we use for sparse smoothing are trained with the noise distribution that is also reported for certification. The localized smoothing models are trained on their minimal noise level, i.e., not with localized noise but with only θ_{\min}^+ and θ_{\min}^- .

Certification We evaluate our certificates on the validation nodes. For all certificates, we use Monte Carlo randomized smoothing (see discussion in § G). We use 1000 samples for making smoothed predictions and $5 \cdot 10^5$ samples for certification. We use the significance parameter α to 0.01, i.e. all certificates hold with probability 0.99. For the naïve isotropic randomized smoothing baseline, we use Holm correction to account for the multiple comparisons problem, which yields strictly better results than Bonferroni correction (see § G.4). For our localized smoothing certificates, we use Bonferroni correction. To parameterize the localized smoothing distribution, we first perform Metis clustering (Karypis & Kumar, 1998) to partition the graph into 5 clusters. We create an affinity ranking by counting the number of edges which are connecting cluster i and j . Specifically, let \mathcal{C} be the set of clusters given by the Metis clustering. Then we count the number of edges between all cluster pairs and denote it by $N_{i,j}$, $i, j \in \mathcal{C}$. If the number of edges of the pair (i, j) is higher than the number for all other pairs $(k, j) \forall j \in \mathcal{C}$, i.e. $N_{i,j} > N_{k,j} \forall k \in \mathcal{C}$, we can say that, due to the homophily assumption, cluster i is the most important one for cluster j . We create this ranking for all pairs and use it to select the noise parameter θ'^- for smoothing the attributes of cluster j while classifying a node of cluster i out of the discrete steps of the linear interpolation between θ_{\min} and θ_{\max} based on its previously defined ranking between the clusters. An example would be, given 11 clusters, $\theta_{\min} = 0.0$, and $\theta_{\max} = 1.0$. If cluster j second most important cluster to i , then we would take the second value out of $\{0.0, 0.1, \dots, 1.0\}$. All results are reported for the relaxed linear programming formulation of our collective certificate (see § E.4). For each cluster, we use $\frac{1}{5}$ of the samples, which corresponds to 200 samples for prediction and 10^5 samples for certification. The collective linear program is solved using MOSEK (version 9.2.46) (MOSEK ApS, 2019) through the CVXPY interface (version 1.1.13) (Diamond & Boyd, 2016).

C.4 HARDWARE AND RUNTIME

The experiments on Pascal-VOC with strictly local models (Fig. 2) were performed using a Xeon E5-2630 v4 CPU @ 2.20GHz, an NVIDIA GTX 1080TI GPU and 128 GB of RAM. All other experiments were performed using an AMD EPYC 7543 CPU @ 2.80GHz, an NVIDIA A100 GPU and 128 GB of RAM.

In all cases, the time needed for obtaining the Monte Carlo samples required by both localized and isotropic smoothing was much larger than the cost of solving the collective linear program.

- For the strictly local model in Fig. 2, taking 153600 samples took 294 s on average. Averaged over all images and adversarial budgets, solving each LP only took 0.91 s.
- For the standard U-Net model in Fig. 4, taking 153600 samples took 70.3 s on average. Each LP took 1.8 s on average.
- For the DeepLabv3 model in Fig. 6b, taking 153600 samples took 1204 s on average. Each LP took 2.78 s on average.
- For the APPNP model in Fig. 5, taking $5 \cdot 10^6$ samples took 1034 s on average. Each LP took 10.9 s on average.

For graphs, the reported time for solving a single instance of the collective linear program is much higher than for image segmentation, even though the graph datasets require fewer variables. That is because we used a different, not as well vectorized formulation of the linear program in CVXPY.

In all cases, the time for calculating the isotropic smoothing certificates and base certificates from the Monte Carlo samples was too small to be measured accurately, since they can be implemented in a few simple vector operations.

D PROOF OF THEOREM 4.2

In the following, we prove Theorem 4.2, i.e. we derive the mixed-integer linear program that underlies our collective certificate and prove that it provides a valid bound on the number of simultaneously robust predictions. The derivation bears some semblance to that of (Schuchardt et al., 2021), in that both use standard techniques to model indicator functions using binary variables and that both convert optimization in input space to optimization in adversarial budget space. Nevertheless, both methods differ in how they encode and evaluate base certificates, ultimately leading to significantly different results (our method encodes each base certificate using only a single linear constraint and does not perform any masking operations).

Theorem 4.2. Given locally smoothed model f , input $\mathbf{x} \in \mathbb{X}^{(D_{\text{in}})}$, smoothed prediction $\mathbf{y} = f(\mathbf{x})$ and base certificates $\mathbb{H}^{(1)}, \dots, \mathbb{H}^{D_{\text{out}}}$ complying with interface Eq. 2, the number of simultaneously robust predictions $\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(\mathbf{x}') = y_n]$ is lower-bounded by

$$\min_{\mathbf{b} \in \mathbb{R}_+^{D_{\text{in}}}, \mathbf{t} \in \{0,1\}^{D_{\text{out}}}} \sum_{n \in \mathbb{T}} t_n \quad (8)$$

$$\text{s.t. } \forall n : \mathbf{b}^T \mathbf{w}^{(n)} \geq (1 - t_n) \eta^{(n)}, \quad \text{sum}\{\mathbf{b}\} \leq \epsilon^p. \quad (9)$$

Proof. We begin by inserting the definition of our perturbation model $\mathbb{B}_{\mathbf{x}}$ and the base certificates $\mathbb{H}^{(n)}$ into Eq. 1.1:

$$\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(\mathbf{x}') = y_n] \geq \min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[\mathbf{x}' \in \mathbb{H}^{(n)}] \quad (10)$$

$$= \min_{\mathbf{x}' \in \mathbb{X}^{D_{\text{in}}}} \sum_{n \in \mathbb{T}} \mathbb{I} \left[\sum_{d=1}^{D_{\text{in}}} w_d^{(n)} \cdot |x'_d - x_d|^p < \eta^{(n)} \right] \text{ s.t. } \sum_{d=1}^{D_{\text{in}}} |x'_d - x_d|^p \leq \epsilon^p. \quad (11)$$

Evidently, input \mathbf{x}' only affects the elementwise distances $|x'_d - x_d|^p$. Rather than optimizing \mathbf{x}' , we can directly optimize these distances, i.e. determine how much adversarial budget is allocated to each input dimension. For this, we define a vector of variables $\mathbf{b} \in \mathbb{R}_+^{D_{\text{in}}}$ (or $\mathbf{b} \in \{0,1\}^{D_{\text{in}}}$ for binary data). Replacing sums with inner products, we can restate Eq. 11 as

$$\min_{\mathbf{b} \in \mathbb{R}_+^{D_{\text{in}}}} \sum_{n \in \mathbb{T}} \mathbb{I}[\mathbf{b}^T \mathbf{w}^{(n)} < \eta^{(n)}] \quad \text{s.t.} \quad \text{sum}\{\mathbf{b}\} \leq \epsilon^p. \quad (12)$$

In a final step, we replace the indicator functions in Eq. 12 with a vector of boolean variables $\mathbf{t} \in \{0,1\}^{D_{\text{out}}}$.

$$\min_{\mathbf{b} \in \mathbb{R}_+^{D_{\text{in}}}, \mathbf{t} \in \{0,1\}^{D_{\text{out}}}} \sum_{n \in \mathbb{T}} t_n \quad (13)$$

$$\text{s.t. } \forall n : \mathbf{b}^T \mathbf{w}^{(n)} \geq (1 - t_n) \eta^{(n)}, \quad \text{sum}\{\mathbf{b}\} \leq \epsilon^p. \quad (14)$$

The first constraint in Eq. 5 ensures that $t_n = 0 \iff \mathbb{I}[\mathbf{b}^T \mathbf{w}^{(n)} \geq \eta^{(n)}]$. Therefore, the optimization problem in Eq. 13 and Eq. 5 is equivalent to Eq. 12, which by transitivity is a lower bound on $\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(\mathbf{x}') = y_n]$. \square

E IMPROVING EFFICIENCY

In this section, we discuss different modifications to our collective certificate that improve its sample efficiency and allow us fine-grained control over the size of the collective linear program. We further discuss a linear relaxation of our collective linear program. All of the modifications preserve the soundness of our collective certificate, i.e. we still obtain a provable bound on the number of predictions that can be simultaneously attacked by an adversary. To avoid constant case distinctions, we first present all results for real-valued data, i.e. $\mathbb{X} = \mathbb{R}$, before mentioning any additional precautions that may be needed when working with binary data.

E.1 SHARING SMOOTHING DISTRIBUTIONS AMONG OUTPUTS

In principle, our proposed certificate allows a different smoothing distribution $\Psi_{\mathbf{x}}^{(n)}$ to be used per output g_n of our base model. In practice, where we have to estimate properties of the smoothed classifier using Monte Carlo methods, this is problematic: Samples cannot be re-used, each of the many outputs requires its own round of sampling. We can increase the efficiency of our localized smoothing approach by partitioning our D_{out} outputs into N_{out} subsets that share the same smoothing distributions. When making smoothed predictions or computing base certificates, we can then reuse the same samples for all outputs within each subsets.

More formally, we partition our D_{out} output dimensions into sets $\mathbb{K}^{(1)}, \dots, \mathbb{K}^{(N_{\text{out}})}$ with

$$\bigcup_{i=1}^{N_{\text{out}}} \mathbb{K}^{(i)} = \{1, \dots, D_{\text{out}}\}. \quad (15)$$

We then associate each set $\mathbb{K}^{(i)}$ with a smoothing distribution $\Psi_{\mathbf{x}}^{(i)}$. For each base model output g_n with $n \in \mathbb{K}^{(i)}$, we then use smoothing distribution $\Psi_{\mathbf{x}}^{(i)}$ to construct the smoothed output f_n , e.g. $f_n(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \Pr_{\mathbf{z} \sim \Psi_{\mathbf{x}}^{(i)}} [f(\mathbf{x} + \mathbf{z}) = y]$ (note that for our variance-constrained certificate we smooth the softmax scores instead, see § 5).

E.2 QUANTIZING CERTIFICATE PARAMETERS

Recall that our base certificates from § 5 are defined by a linear inequality: A prediction $y_n = f_n(\mathbf{x})$ is robust to a perturbed input $\mathbf{x}' \in \mathbb{X}^{D_{\text{in}}}$ if $\sum_{d=1}^D w_d^{(n)} \cdot |x'_d - x_d|^p < \eta^{(n)}$, for some $p \geq 0$. The weight vectors $\mathbf{w}^{(n)} \in \mathbb{R}^{D_{\text{in}}}$ only depend on the smoothing distributions. A side effect of sharing the same distribution $\Psi_{\mathbf{x}}^{(i)}$ among all outputs from a set $\mathbb{K}^{(i)}$, as discussed in the previous section, is that the outputs also share the same weight vector $\mathbf{w}^{(i)} \in \mathbb{R}^{D_{\text{in}}}$ with $\forall n \in \mathbb{K}^{(i)} : \mathbf{w}^{(i)} = \mathbf{w}^{(n)}$. Thus, for all smoothed outputs f_n with $n \in \mathbb{K}^{(i)}$, the smoothed prediction y_n is robust if $\sum_{d=1}^D w_d^{(i)} \cdot |x'_d - x_d|^p < \eta^{(n)}$.

Evidently, the base certificates for outputs from a set $\mathbb{K}^{(i)}$ only differ in their parameter $\eta^{(n)}$. Recall that in our collective linear program we use a vector of variables $\mathbf{t} \in \{0, 1\}^{D_{\text{out}}}$ to indicate which predictions are robust according to their base certificates (see Theorem 4.2). If there are two outputs f_n and f_m with $\eta^{(n)} = \eta^{(m)}$, then f_n and f_m have the same base certificate and their robustness can be modelled by the same indicator variable. Conversely, for each set of outputs $\mathbb{K}^{(i)}$, we only need one indicator variable per unique $\eta^{(n)}$. By quantizing the $\eta^{(n)}$ within each subset $\mathbb{K}^{(i)}$ (for example by defining equally sized bins between $\min_{n \in \mathbb{K}^{(i)}} \eta^{(n)}$ and $\max_{n \in \mathbb{K}^{(i)}} \eta^{(n)}$), we can ensure that there is always a fixed number N_{bins} of indicator variables per subset. This way, we can reduce the number of indicator variables from D_{out} to $N_{\text{out}} \cdot N_{\text{bins}}$.

To implement this idea, we define a matrix of thresholds $\mathbf{E} \in \mathbb{R}^{N_{\text{out}} \times N_{\text{bins}}}$ with $\forall i : \min \{\mathbf{E}_{i,:}\} \leq \min_{n \in \mathbb{K}^{(i)}} (\{\eta^{(n)} \mid n \in \mathbb{K}^{(i)}\})$. We then define a function $\xi : \{1, \dots, N_{\text{out}}\} \times \mathbb{R} \rightarrow \mathbb{R}$ with

$$\xi(i, \eta) = \max (\{ \mathbf{E}_{i,j} \mid j \in \{1, \dots, N_{\text{bins}} \wedge \mathbf{E}_{i,j} < \eta \}) \quad (16)$$

that quantizes base certificate parameter η from output subset $\mathbb{K}^{(i)}$ by mapping it to the next smallest threshold in $\mathbf{E}_{i,:}$. We can then bound the collective robustness of the targeted dimensions \mathbb{T} of our

prediction vector $\mathbf{y} = f(\mathbf{x})$ as follows:

$$\min \sum_{i \in \{1, \dots, N_{\text{out}}\}} \sum_{j \in \{1, \dots, N_{\text{bins}}\}} T_{i,j} \left| \left\{ n \in \mathbb{T} \cap \mathbb{K}^{(i)} \mid \xi(i, \eta^{(n)}) = E_{i,j} \right\} \right| \quad (17)$$

$$\text{s.t. } \forall i, j : \mathbf{b}^T \mathbf{w}^{(i)} \geq (1 - T_{i,j}) E_{i,j}, \quad \text{sum}\{\mathbf{b}\} \leq \epsilon^p \quad (18)$$

$$\mathbf{b} \in \mathbb{R}_+^{D_{\text{in}}}, \quad \mathbf{T} \in \{0, 1\}^{N_{\text{out}} \times N_{\text{bins}}}. \quad (19)$$

Constraint Eq. 18 ensures that $T_{i,j}$ is only set to 0 if $\mathbf{b}^T \mathbf{w}^{(i)} \geq E_{i,j}$, i.e. all predictions from subset $\mathbb{K}^{(i)}$ whose base certificate parameter $\eta^{(n)}$ is quantized to $E_{i,j}$ are no longer robust. When this is the case, the objective function decreases by the number of these predictions. For $N_{\text{out}} = D_{\text{out}}$, $N_{\text{bins}} = 1$ and $E_{n,1} = \eta^{(n)}$, we recover our general certificate from Theorem 4.2. Note that, if the quantization maps any parameter $\eta^{(n)}$ to a smaller number, the base certificate $\mathbb{H}^{(n)}$ becomes more restrictive, i.e. y_n is considered robust to a smaller set of perturbed inputs. Thus, Eq. 17 is a lower bound on our general certificate from Theorem 4.2.

E.3 SHARING NOISE LEVELS AMONG INPUTS

Similar to how partitioning the output dimensions allows us to control the number of output variables \mathbf{t} , partitioning the input dimensions and using the same noise level within each partition allows us to control the number of budget variables \mathbf{b} .

Assume that we have partitioned our output dimensions into N_{out} subsets $\mathbb{K}^{(1)}, \dots, \mathbb{K}^{(N_{\text{out}})}$, with outputs in each subset sharing the same smoothing distribution $\Psi_{\mathbf{x}}^{(i)}$, as explained in § E.1. Let us now define N_{in} input subsets $\mathbb{J}^{(1)}, \dots, \mathbb{J}^{(N_{\text{in}})}$ with

$$\dot{\bigcup}_{l=1}^{N_{\text{in}}} \mathbb{J}^{(l)} = \{1, \dots, D_{\text{out}}\}. \quad (20)$$

Recall that a prediction $y_n = f_n(\mathbf{x})$ with $n \in \mathbb{K}^{(i)}$ is robust to a perturbed input $\mathbf{x}' \in \mathbb{X}^{D_{\text{in}}}$ if $\sum_{d=1}^{D_{\text{in}}} w_d^{(i)} \cdot |x'_d - x_d|^p < \eta^{(n)}$ and that the weight vectors $\mathbf{w}^{(i)}$ only depend on the smoothing distributions. Assume that we choose each smoothing distribution $\Psi_{\mathbf{x}}^{(i)}$ such that $\forall l \in \{1, \dots, N_{\text{in}}\}, \forall d, d' \in \mathbb{J}^{(l)} : w_d^{(i)} = w_{d'}^{(i)}$, i.e. all input dimensions within each set $\mathbb{J}^{(l)}$ have the same weight. This can be achieved by choosing $\Psi_{\mathbf{x}}^{(i)}$ so that all dimensions in each input subset $\mathbb{J}^{(l)}$ are smoothed with the noise level (note that we can still use a different smoothing distribution $\Psi_{\mathbf{x}}^{(i)}$ for each set of outputs $\mathbb{K}^{(i)}$). For example, one could use a Gaussian distribution with covariance matrix $\Sigma = \text{diag}(\boldsymbol{\sigma})^2$ with $\forall l \in \{1, \dots, N_{\text{in}}\}, \forall d, d' \in \mathbb{J}^{(l)} : \sigma_d = \sigma_{d'}$.

In this case, the evaluation of our base certificates can be simplified. Prediction $y_n = f_n(\mathbf{x})$ with $n \in \mathbb{K}^{(n)}$ is robust to a perturbed input $\mathbf{x}' \in \mathbb{X}^{D_{\text{in}}}$ if

$$\sum_{d=1}^{D_{\text{in}}} w_d^{(i)} \cdot |x'_d - x_d|^p < \eta^{(n)} \quad (21)$$

$$= \sum_{l=1}^{N_{\text{in}}} \left(u^{(i)} \cdot \sum_{d \in \mathbb{J}^{(l)}} |x'_d - x_d|^p \right) < \eta^{(n)}, \quad (22)$$

with $\mathbf{u} \in \mathbb{R}_+^{N_{\text{in}}}$ and $\forall i \in \{1, \dots, N_{\text{out}}\}, \forall l \in \{1, \dots, N_{\text{in}}\}, \forall d \in \mathbb{J}^{(l)} : u_l^i = w_d^i$. That is, we can replace each weight vector $\mathbf{w}^{(i)}$ that has one weight $w_d^{(i)}$ per input dimension d with a smaller weight vector $\mathbf{u}^{(i)}$ featuring one weight $u_l^{(i)}$ per input subset $\mathbb{J}^{(l)}$.

For our linear program, this means that we no longer need a budget vector $\mathbf{b} \in \mathbb{R}_+^{D_{\text{in}}}$ to model the elementwise distance $|x'_d - x_d|^p$ in each dimension d . Instead, we can use a smaller budget vector $\mathbf{b} \in \mathbb{R}_+^{N_{\text{in}}}$ to model the overall distance within each input subset $\mathbb{J}^{(l)}$, i.e. $b^{(l)} = \sum_{d \in \mathbb{J}^{(l)}} |x'_d - x_d|^p$. Combined with the quantization of certificate parameters from the previous section, our optimization

problem becomes

$$\min \sum_{i \in \{1, \dots, N_{\text{out}}\}} \sum_{j \in \{1, \dots, N_{\text{bins}}\}} T_{i,j} \left| \left\{ n \in \mathbb{T} \cap \mathbb{K}^{(i)} \mid \xi(i, \eta^{(n)}) = E_{i,j} \right\} \right| \quad (23)$$

$$\text{s.t. } \forall i, j : \mathbf{b}^T \mathbf{u}^{(i)} \geq (1 - T_{i,j}) E_{i,j}, \quad \text{sum}\{\mathbf{b}\} \leq \epsilon^p, \quad (24)$$

$$\mathbf{b} \in \mathbb{R}_+^{N_{\text{in}}}, \quad \mathbf{T} \in \{0, 1\}^{N_{\text{out}} \times N_{\text{bins}}}. \quad (25)$$

with $\mathbf{u} \in \mathbb{R}^{N_{\text{in}}}$ and $\forall i \in \{1, \dots, N_{\text{out}}\}, \forall l \in \{1, \dots, N_{\text{in}}\}, \forall d \in \mathbb{J} : u_l^i = w_d^i$. For $N_{\text{out}} = D_{\text{out}}, N_{\text{in}} = D_{\text{in}}, N_{\text{bins}} = 1$ and $E_{n,1} = \eta^{(n)}$, we recover our general certificate from Theorem 4.2.

When certifying robustness for binary data, we impose different constraints on \mathbf{b} . To model that the adversary can not flip more bits than are present within each subset, we use a budget vector $\mathbf{b} \in \mathbb{N}_0^{N_{\text{in}}}$ with $\forall l \in \{1, \dots, N_{\text{in}}\} : b_l \leq |\mathbb{J}^{(l)}|$, instead of a continuous budget vector $\mathbf{b} \in \mathbb{R}_+^{N_{\text{in}}}$.

E.4 LINEAR RELAXATION

Combining the previous steps allows us to reduce the number of problem variables and linear constraints from $D_{\text{in}} + D_{\text{out}}$ and $D_{\text{out}} + 1$ to $N_{\text{in}} + N_{\text{out}} \cdot N_{\text{bins}}$ and $N_{\text{out}} \cdot N_{\text{bins}} + 1$, respectively. Still, finding an optimal solution to the mixed-integer linear program may be too expensive. One can obtain a lower bound on the optimal value and thus a valid, albeit more pessimistic, robustness certificate by relaxing all discrete variables to be continuous.

When using the general certificate from Theorem 4.2, the binary vector $\mathbf{t} \in \{0, 1\}^{D_{\text{out}}}$ can be relaxed to $\mathbf{t} \in [0, 1]^{D_{\text{out}}}$. When using the certificate with quantized base certificate parameters from § E.2 or § E.3, the binary matrix $\mathbf{T} \in \{0, 1\}^{N_{\text{out}} \times N_{\text{bins}}}$ can be relaxed to $\mathbf{T} \in [0, 1]^{N_{\text{out}} \times N_{\text{bins}}}$. Conceptually, this means that predictions can be partially certified, i.e. $t_n \in (0, 1)$ or $T_{i,j} \in (0, 1)$. In particular, a prediction can be partially certified even if we know that is impossible to attack under the collective perturbation model $\mathbb{B}_{\mathbf{x}} = \{\mathbf{x}' \in \mathbb{X}^{D_{\text{in}}} \mid \|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon\}$. Just like Schuchardt et al. (2021), who encountered the same problem with their collective certificate, we circumvent this issue by first computing a set $\mathbb{L} \subseteq \mathbb{T}$ of all targeted predictions in \mathbb{T} that are guaranteed to always be robust under the collective perturbation model:

$$\mathbb{L} = \left\{ n \in \mathbb{T} \mid \left(\max_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{d=1}^D w_d^{(n)} \cdot |x'_d - x_d|^p \right) < \eta^{(n)} \right\} \quad (26)$$

$$= \left\{ n \in \mathbb{T} \mid \max_n \{ \mathbf{w}^{(n)} \} \cdot \epsilon^p < \eta^{(n)} \right\}. \quad (27)$$

The equality follows from the fact that the most effective way of attacking a prediction is to allocate all adversarial budget to the least robust dimension, i.e. the dimension with the largest weight. Because we know that all predictions with indices in \mathbb{L} are robust, we do not have to include them in the collective optimization problem and can instead compute

$$|\mathbb{L}| + \min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T} \setminus \mathbb{L}} \mathbb{I}[\mathbf{x}' \in \mathbb{H}^{(n)}]. \quad (28)$$

The r.h.s. optimization can be solved using the general collective certificate from Theorem 4.2 or any of the more efficient, modified certificates from previous sections.

When using the general collective certificate from Theorem 4.2 with binary data, the budget variables $\mathbf{b} \in \{0, 1\}^{D_{\text{in}}}$ can be relaxed to $\mathbf{b} \in [0, 1]^{D_{\text{in}}}$. When using the modified collective certificate from § E.3, the budget variables with $\mathbf{b} \in \mathbb{N}_0^{N_{\text{in}}}$ can be relaxed to $\mathbf{b} \in \mathbb{R}_+^{N_{\text{in}}}$. The additional constraint $\forall l \in \{1, \dots, N_{\text{in}}\} : b_l \leq |\mathbb{J}^{(l)}|$ can be kept in order to model that the adversary cannot flip (or partially flip) more bits than are present within each input subset $\mathbb{J}^{(l)}$.

F BASE CERTIFICATES

In the following, we show why the base certificates discussed in § 5 and summarized in Table 1 hold. In § F.3.2 we further present a base certificate (and corresponding collective certificate) that can distinguish between adversarial addition and deletion of bits in binary data.

F.1 GAUSSIAN SMOOTHING FOR l_2 PERTURBATIONS OF CONTINUOUS DATA

Proposition F.1. *Given an output $g_n : \mathbb{R}^{D_{\text{in}}} \rightarrow \mathbb{Y}$, let $f_n(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{x}, \Sigma)} [g_n(\mathbf{z}) = y]$ be the corresponding smoothed output with $\Sigma = \operatorname{diag}(\boldsymbol{\sigma})^2$ and $\boldsymbol{\sigma} \in \mathbb{R}_+^{D_{\text{in}}}$. Given an input $\mathbf{x} \in \mathbb{R}^{D_{\text{in}}}$ and smoothed prediction $y_n = f_n(\mathbf{x})$, let $q = \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{x}, \Sigma)} [g_n(\mathbf{z}) = y_n]$. Then, $\forall \mathbf{x}' \in \mathbb{H}^{(n)} : f_n(\mathbf{x}') = y_n$ with $\mathbb{H}^{(n)}$ defined as in Eq. 2, $w_d = \frac{1}{\sigma_d^2}$, $\eta = (\Phi^{-1}(q))^2$ and $p = 2$.*

Proof. Based on the definition of the base certificate interface, we need to show that, $\forall \mathbf{x}' \in \mathbb{H} : f_n(\mathbf{x}') = y_n$ with

$$\mathbb{H} = \left\{ \mathbf{x}' \in \mathbb{R}^{D_{\text{in}}} \mid \sum_{d=1}^{D_{\text{in}}} \frac{1}{\sigma_d^2} \cdot |x_d - x'_d|^2 < (\Phi^{-1}(q))^2 \right\}. \quad (29)$$

Eiras et al. (2022) have shown that under the same conditions as above, but with a general covariance matrix $\Sigma \in \mathbb{R}_+^{D_{\text{in}} \times D_{\text{in}}}$, a prediction y_n is certifiably robust to a perturbed input \mathbf{x}' if

$$\sqrt{(\mathbf{x} - \mathbf{x}') \Sigma^{-1} (\mathbf{x} - \mathbf{x}')} < \frac{1}{2} (\Phi^{-1}(q) - \Phi^{-1}(q')), \quad (30)$$

where $q' = \max_{y'_n \neq y_n} \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{x}, \Sigma)} [g_n(\mathbf{z}) = y'_n]$ is the probability of the second most likely prediction under the smoothing distribution. Because the probabilities of all possible predictions have to sum up to 1, we have $q' \leq 1 - q$. Since Φ^{-1} is monotonically increasing, we can obtain a lower bound on the r.h.s. of Eq. 30 and thus a more pessimistic certificate by substituting $1 - q$ for q' (deriving such a "binary certificate" from a "multiclass certificate" is common in randomized smoothing and was already discussed in (Cohen et al., 2019)):

$$\sqrt{(\mathbf{x} - \mathbf{x}') \Sigma^{-1} (\mathbf{x} - \mathbf{x}')} < \frac{1}{2} (\Phi^{-1}(q) - \Phi^{-1}(1 - q)), \quad (31)$$

In our case, Σ is a diagonal matrix $\operatorname{diag}(\boldsymbol{\sigma})^2$ with $\boldsymbol{\sigma} \in \mathbb{R}_+^{D_{\text{in}}}$. Thus Eq. 31 is equivalent to

$$\sqrt{\sum_{d=1}^{D_{\text{in}}} (x_d - x'_d) \frac{1}{\sigma_d^2} (x_d - x'_d)} < \frac{1}{2} (\Phi^{-1}(q) - \Phi^{-1}(1 - q)). \quad (32)$$

Finally, using the fact that $\Phi^{-1}(q) - \Phi^{-1}(1 - q) = 2\Phi^{-1}(q)$ and eliminating the square root shows that we are certifiably robust if

$$\sum_{d=1}^{D_{\text{in}}} \frac{1}{\sigma_d^2} \cdot |x_d - x'_d|^2 < (\Phi^{-1}(q))^2. \quad (33)$$

□

Table 1: Base certificates complying with interface Eq. 2 with parameters $w^{(n)}$ and $\eta^{(n)}$. Here, $y_n = f_n(\mathbf{x})$ is the prediction of $f_n(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} q_{n,y}$. With the l_0 certificate, $g_n(\mathbf{z})_y$ refers to the softmax score of class y and $\zeta = \operatorname{Var}_{\mathbf{z} \sim \mathcal{F}(\mathbf{x}, \boldsymbol{\theta})} [g_n(\mathbf{z})_{y_n}]$ is the variance of y_n 's softmax score.

Norm	$\Psi_{\mathbf{x}}^{(n)}$	$q_{n,y}$	$w_d^{(n)}$	$\eta^{(n)}$
l_2	$\mathcal{N}(\mathbf{x}, \operatorname{diag}(\mathbf{s})^2)$	$\Pr_{\mathbf{z} \sim \Psi_{\mathbf{x}}^{(n)}} [g_n(\mathbf{z}) = y]$	$\frac{1}{s_d^2}$	$(\Phi^{-1}(q_{n,y_n}))^2$
l_1	$\mathcal{U}(\mathbf{x}, \boldsymbol{\lambda})$	$\Pr_{\mathbf{z} \sim \Psi_{\mathbf{x}}^{(n)}} [g_n(\mathbf{z}) = y]$	$\frac{1}{\lambda_d}$	$\Phi^{-1}(q_{n,y_n})$
l_0	$\mathcal{F}(\mathbf{x}, \boldsymbol{\theta})$	$\mathbb{E}_{\mathbf{z} \sim \Psi_{\mathbf{x}}^{(n)}} [g_n(\mathbf{z})_y]$	$\ln \left(\frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d} \right)$	$\ln \left(1 + \frac{1}{\zeta} (q_{n,y_n} - \frac{1}{2})^2 \right)$

F.2 UNIFORM SMOOTHING FOR l_1 PERTURBATIONS OF CONTINUOUS DATA

An alternative base certificate for l_1 perturbations is again due to Eiras et al. (2022). Using uniform instead of Gaussian noise allows us to collectively certify robustness to l_1 -norm-bound perturbations. In the following $\mathcal{U}(\mathbf{x}, \boldsymbol{\lambda})$ with $\mathbf{x} \in \mathbb{R}^D$, $\boldsymbol{\lambda} \in \mathbb{R}_+^D$ refers to a vector-valued random distribution in which the d -th element is uniformly distributed in $[x_d - \lambda_d, x_d + \lambda_d]$.

Proposition F.2. *Given an output $g_n : \mathbb{R}^{D_{\text{in}}} \rightarrow \mathbb{Y}$, let $f(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \Pr_{\mathbf{z} \sim \mathcal{U}(\mathbf{x}, \boldsymbol{\lambda})} [g(\mathbf{z}) = y]$ be the corresponding smoothed classifier with $\boldsymbol{\lambda} \in \mathbb{R}_+^{D_{\text{in}}}$. Given an input $\mathbf{x} \in \mathbb{R}^{D_{\text{in}}}$ and smoothed prediction $y = f(\mathbf{x})$, let $p = \Pr_{\mathbf{z} \sim \mathcal{U}(\mathbf{x}, \boldsymbol{\lambda})} [g(\mathbf{z}) = y]$. Then, $\forall \mathbf{x}' \in \mathbb{H}^{(n)} : f_n(\mathbf{x}') = y_n$ with $\mathbb{H}^{(n)}$ defined as in Eq. 2, $w_d = 1/\lambda_d$, $\eta = \Phi^{-1}(q)$ and $p = 1$.*

Proof. Based on the definition of $\mathbb{H}^{(n)}$, we need to prove that $\forall \mathbf{x}' \in \mathbb{H} : f_n(\mathbf{x}') = y_n$ with

$$\mathbb{H} = \left\{ \mathbf{x}' \in \mathbb{R}^{D_{\text{in}}} \mid \sum_{d=1}^{D_{\text{in}}} \frac{1}{\lambda_d} \cdot |x_d - x'_d| < \Phi^{-1}(q) \right\}, \quad (34)$$

Eiras et al. (2022) have shown that under the same conditions as above, a prediction y_n is certifiably robust to a perturbed input \mathbf{x}' if

$$\sum_{d=1}^{D_{\text{in}}} \frac{1}{\lambda_d} \cdot (x_d - x'_d) < \frac{1}{2} (\Phi^{-1}(q) - \Phi^{-1}(1-q)), \quad (35)$$

where $q' = \max_{y'_n \neq y_n} \Pr_{\mathbf{z} \sim \mathcal{U}(\mathbf{x}, \boldsymbol{\lambda})} [g_n(\mathbf{z}) = y'_n]$ is the probability of the second most likely prediction under the smoothing distribution. As in our previous proof for Gaussian smoothing, we can obtain a more pessimistic certificate by substituting $1-q$ for q' . Since $\Phi^{-1}(q) - \Phi^{-1}(1-q) = 2\Phi^{-1}(q)$ and all λ_d are non-negative, we know that our prediction is certifiably robust if

$$\sum_{d=1}^{D_{\text{in}}} \frac{1}{\lambda_d} \cdot |x_d - x'_d| < \Phi^{-1}(p). \quad (36)$$

□

F.3 VARIANCE-CONSTRAINED CERTIFICATION

In the following, we derive the general variance-constrained randomized smoothing certificate from Theorem 5.1, before discussing specific certificates for binary data in § F.3.1 and § F.3.2.

Variance smoothing assumes that we make predictions by randomly smoothing a base model's softmax scores. That is, given base model $g : \mathbb{X} \rightarrow \Delta_{|\mathbb{Y}|}$ mapping from an arbitrary discrete input space \mathbb{X} to scores from the $(|\mathbb{Y}| - 1)$ -dimensional probability simplex $\Delta_{|\mathbb{Y}|}$, we define the smoothed classifier $f(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [g(\mathbf{z})_y]$. Here, $\Psi(\mathbf{x})$ is an arbitrary distribution over \mathbb{X} parameterized by \mathbf{x} , e.g a Normal distribution with mean \mathbf{x} . The smoothed classifier does not return the most likely prediction, but the prediction associated with the highest expected softmax score.

Given an input $\mathbf{x} \in \mathbb{X}$, smoothed prediction $y = f(\mathbf{x})$ and a perturbed input $\mathbf{x}' \in \mathbb{X}$, we want to determine whether $f(\mathbf{x}') = y$. By definition of our smoothed classifier, we know that $f(\mathbf{x}') = y$ if y is the label with the highest expected softmax score. In particular, we know that $f(\mathbf{x}') = y$ if y 's softmax score is larger than all other softmax scores combined, i.e.

$$\mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x}')} [g(\mathbf{z})_y] > 0.5 \implies f(\mathbf{x}') = y. \quad (37)$$

Computing $\mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x}')} [g(\mathbf{z})_y]$ exactly is usually not tractable – especially if we later want to evaluate robustness to many \mathbf{x}' from a whole perturbation model $\mathbb{B} \subseteq \mathbb{X}$. Therefore, we compute a lower bound on $\mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x}')} [g(\mathbf{z})_y]$. If even this lower bound is larger than 0.5, we know that prediction y is certainly robust. For this, we define a set of functions \mathbb{F} with $g_y \in \mathbb{H}$ and compute the minimum softmax score across all functions from \mathbb{F} :

$$\min_{h \in \mathbb{F}} \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x}')} [h(\mathbf{z})] > 0.5 \implies f(\mathbf{x}') = y. \quad (38)$$

For our variance smoothing approach, we define \mathbb{F} to be the set of all functions that have a larger or equal expected value and a smaller or equal variance under $\Psi(\mathbf{x})$, compared to our base model g . Let $\mu = \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [g(\mathbf{z})_y]$ be the expected softmax score of our base model g for label y . Let $\zeta = \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} \left[(g(\mathbf{z})_y - \nu)^2 \right]$ be the expected squared distance of the softmax score from a scalar $\nu \in \mathbb{R}$. (Choosing $\nu = \mu$ yields the variance of the softmax score. An arbitrary ν is only needed for technical reasons related to Monte Carlo estimation § G.2). Then, we define

$$\mathbb{F} = \left\{ h : \mathbb{X} \rightarrow \mathbb{R} \mid \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [h(\mathbf{z})] \geq \mu \wedge \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} \left[(h(\mathbf{z}) - \nu)^2 \right] \leq \zeta \right\} \quad (39)$$

Clearly, by the definition of μ and ζ , we have $g_y \in \mathbb{F}$. Note that we do not restrict functions from \mathbb{H} to the domain $[0, 1]$, but allow arbitrary real-valued outputs.

By evaluating Eq. 37 with \mathbb{F} defined as in Eq. 38, we can determine if our prediction is robust. To compute the optimal value, we need the following two Lemmata:

Lemma F.3. *Given a discrete set \mathbb{X} and the set Π of all probability mass functions over \mathbb{X} , any two probability mass functions $\pi_1, \pi_2 \in \Pi$ fulfill*

$$\sum_{z \in \mathbb{X}} \frac{\pi_2(z)}{\pi_1(z)} \pi_2(z) \geq 1. \quad (40)$$

Proof. For a fixed probability mass function π_1 , Eq. 40 is lower-bounded by the minimal expected likelihood ratio that can be achieved by another $\tilde{\pi}(z) \in \Pi$:

$$\sum_{z \in \mathbb{X}} \frac{\pi_2(z)}{\pi_1(z)} \pi_2(z) \geq \min_{\tilde{\pi} \in \Pi} \sum_{z \in \mathbb{X}} \frac{\tilde{\pi}(z)}{\pi_1(z)} \tilde{\pi}(z). \quad (41)$$

The r.h.s. term can be expressed as the constrained optimization problem

$$\min_{\tilde{\pi}} \sum_{z \in \mathbb{X}} \frac{\tilde{\pi}(z)}{\pi_1(z)} \tilde{\pi}(z) \quad \text{s.t.} \quad \sum_{z \in \mathbb{X}} \tilde{\pi}(z) = 1 \quad (42)$$

with the corresponding dual problem

$$\max_{\lambda \in \mathbb{R}} \min_{\tilde{\pi}} \sum_{z \in \mathbb{X}} \frac{\tilde{\pi}(z)}{\pi_1(z)} \tilde{\pi}(z) + \lambda \left(-1 + \sum_{z \in \mathbb{X}} \tilde{\pi}(z) \right). \quad (43)$$

The inner problem is convex in each $\tilde{\pi}(z)$. Taking the gradient w.r.t. to $\tilde{\pi}(z)$ for all $z \in \mathbb{X}$ shows that it has its minimum at $\forall z \in \mathbb{X} : \tilde{\pi}(z) = -\frac{\lambda \pi_1(z)}{2}$. Substituting into Eq. 43 results in

$$\max_{\lambda \in \mathbb{R}} \sum_{z \in \mathbb{X}} \frac{\lambda^2 \pi_1(z)^2}{4 \pi_1(z)} + \lambda \left(-1 - \sum_{z \in \mathbb{X}} \frac{\lambda \pi_1(z)}{2} \right) \quad (44)$$

$$= \max_{\lambda \in \mathbb{R}} -\lambda^2 \sum_{z \in \mathbb{X}} \frac{\pi_1(z)}{4} - \lambda \quad (45)$$

$$= \max_{\lambda \in \mathbb{R}} -\frac{\lambda^2}{4} - \lambda \quad (46)$$

$$= 1. \quad (47)$$

Eq. 46 follows from the fact that $\pi_1(z)$ is a valid probability mass function. Due to duality, the optimal dual value 1 is a lower bound on the optimal value of our primal problem Eq. 40. \square

Lemma F.4. *Given a probability distribution \mathcal{D} over a \mathbb{R} and a scalar $\nu \in \mathbb{R}$, let $\mu = \mathbb{E}_{z \sim \mathcal{D}} [z]$ and $\xi = \mathbb{E}_{z \sim \mathcal{D}} \left[(z - \nu)^2 \right]$. Then $\xi \geq (\mu - \nu)^2$*

Proof. Using the definitions of μ and ξ , as well as some simple algebra, we can show:

$$\xi \geq (\mu - \nu)^2 \quad (48)$$

$$\iff \mathbb{E}_{z \sim \mathcal{D}} \left[(z - \nu)^2 \right] \geq \mu^2 - 2\mu\nu + \nu^2 \quad (49)$$

$$\iff \mathbb{E}_{z \sim \mathcal{D}} \left[z^2 - 2z\nu + \nu^2 \right] \geq \mu^2 - 2\mu\nu + \nu^2 \quad (50)$$

$$\iff \mathbb{E}_{z \sim \mathcal{D}} \left[z^2 - 2z\nu + \nu^2 \right] \geq \mu^2 - 2\mu\nu + \nu^2 \quad (51)$$

$$\iff \mathbb{E}_{z \sim \mathcal{D}} \left[z^2 \right] - 2\mu\nu + \nu^2 \geq \mu^2 - 2\mu\nu + \nu^2 \quad (52)$$

$$\iff \mathbb{E}_{z \sim \mathcal{D}} \left[z^2 \right] \geq \mu^2 \quad (53)$$

It is well known for the variance that $\mathbb{E}_{z \sim \mathcal{D}} \left[(z - \mu)^2 \right] = \mathbb{E}_{z \sim \mathcal{D}} \left[z^2 \right] - \mu^2$. Because the variance is always non-negative, the above inequality holds. \square

Using the previously described approach and lemmata, we can show the soundness of the following robustness certificate:

Theorem 5.1 (Variance-constrained certification). a function $g : \mathbb{X} \rightarrow \Delta_{|\mathbb{Y}|}$ mapping from discrete set \mathbb{X} to scores from the $(|\mathbb{Y}| - 1)$ -dimensional probability simplex, let $f(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \mathbb{E}_{z \sim \Psi_{\mathbf{x}}} [g(z)_y]$ with smoothing distribution $\Psi_{\mathbf{x}}$ and probability mass function $\pi_{\mathbf{x}}(\mathbf{z}) = \Pr_{\tilde{z} \sim \Psi_{\mathbf{x}}} [\tilde{z} = \mathbf{z}]$. Given an input $\mathbf{x} \in \mathbb{X}$ and smoothed prediction $y = f(\mathbf{x})$, let $\mu = \mathbb{E}_{z \sim \Psi_{\mathbf{x}}} [g(z)_y]$ and $\zeta = \mathbb{E}_{z \sim \Psi_{\mathbf{x}}} \left[(g(z)_y - \nu)^2 \right]$ with $\nu \in \mathbb{R}$. Assuming $\nu \leq \mu$, then $f(\mathbf{x}') = y$ if

$$\sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})}{\pi_{\mathbf{x}}(\mathbf{z})} \cdot \pi_{\mathbf{x}'}(\mathbf{z}) < 1 + \frac{1}{\zeta - (\mu - \nu)^2} \left(\mu - \frac{1}{2} \right). \quad (54)$$

Proof. Following our discussion above, we know that $f(\mathbf{x}') = y$ if $\mathbb{E}_{z \sim \Psi(\mathbf{x}')} [g(z)_y] > 0.5$ with \mathbb{F} defined as in Eq. 39. We can compute a (tight) lower bound on $\min_{h \in \mathbb{F}} \mathbb{E}_{z \sim \Psi(\mathbf{x}')} [h(\mathbf{z})]$ by following the functional optimization approach for randomized smoothing proposed by Zhang et al. (2020). That is, we solve a dual problem in which we optimize the value $h(\mathbf{z})$ for each $\mathbf{z} \in \mathbb{X}$. By the definition of the set \mathbb{F} , our optimization problem is

$$\min_{h: \mathbb{X} \rightarrow \mathbb{R}} \mathbb{E}_{z \sim \Psi(\mathbf{x}')} [h(\mathbf{z})] \quad (55)$$

$$\text{s.t. } \mathbb{E}_{z \sim \Psi(\mathbf{x})} [h(\mathbf{z})] \geq \mu, \quad \mathbb{E}_{z \sim \Psi(\mathbf{x})} \left[(h(\mathbf{z}) - \nu)^2 \right] \leq \zeta. \quad (56)$$

The corresponding dual problem with dual variables $\alpha, \beta \geq 0$ is

$$\begin{aligned} & \max_{\alpha, \beta \geq 0} \min_{h: \mathbb{X} \rightarrow \mathbb{R}} \mathbb{E}_{z \sim \Psi(\mathbf{x}')} [h(\mathbf{z})] \\ & + \alpha \left(\mu - \mathbb{E}_{z \sim \Psi(\mathbf{x})} [h(\mathbf{z})] \right) + \beta \left(\mathbb{E}_{z \sim \Psi(\mathbf{x})} \left[(h(\mathbf{z}) - \nu)^2 \right] - \zeta \right). \end{aligned} \quad (57)$$

We first move all terms that don't involve h out of the inner optimization problem:

$$= \max_{\alpha, \beta \geq 0} \alpha\mu - \beta\zeta + \min_{h: \mathbb{X} \rightarrow \mathbb{R}} \mathbb{E}_{z \sim \Psi(\mathbf{x}')} [h(\mathbf{z})] - \alpha \mathbb{E}_{z \sim \Psi(\mathbf{x})} [h(\mathbf{z})] + \beta \mathbb{E}_{z \sim \Psi(\mathbf{x})} \left[(h(\mathbf{z}) - \nu)^2 \right]. \quad (58)$$

Writing out the expectation terms and combining them into one sum (or – in the case of continuous \mathbb{X} – one integral), our dual problem becomes

$$= \max_{\alpha, \beta \geq 0} \alpha\mu - \beta\zeta + \min_{h: \mathbb{X} \rightarrow \mathbb{R}} \sum_{\mathbf{z} \in \mathbb{X}} h(\mathbf{z})\pi_{\mathbf{x}'}(\mathbf{z}) - \alpha h(\mathbf{z})\pi_{\mathbf{x}}(\mathbf{z}) + \beta (h(\mathbf{z}) - \nu)^2 \pi_{\mathbf{x}}(\mathbf{z}) \quad (59)$$

(recall that $\pi_{\mathbf{x}'}$ and $\pi_{\mathbf{x}}$ refer to the probability mass functions of the smoothing distributions). The inner optimization problem can be solved by finding the optimal $h(\mathbf{z})$ in each point \mathbf{z} :

$$= \max_{\alpha, \beta \geq 0} \alpha\mu - \beta\zeta + \sum_{\mathbf{z} \in \mathbb{X}} \min_{h(\mathbf{z}) \in \mathbb{R}} h(\mathbf{z})\pi_{\mathbf{x}'}(\mathbf{z}) - \alpha h(\mathbf{z})\pi_{\mathbf{x}}(\mathbf{z}) + \beta (h(\mathbf{z}) - \nu)^2 \pi_{\mathbf{x}}(\mathbf{z}). \quad (60)$$

Because $\beta \geq 0$, each inner optimization problem is convex in $h(\mathbf{z})$. We can thus find the optimal $h^*(\mathbf{z})$ by setting the derivative to zero:

$$\frac{d}{dh(\mathbf{z})} h(\mathbf{z})\pi_{\mathbf{x}'}(\mathbf{z}) - \alpha h(\mathbf{z})\pi_{\mathbf{x}}(\mathbf{z}) + \beta (h(\mathbf{z}) - \nu)^2 \pi_{\mathbf{x}}(\mathbf{z}) \stackrel{!}{=} 0 \quad (61)$$

$$\Leftrightarrow \pi_{\mathbf{x}'}(\mathbf{z}) - \alpha \pi_{\mathbf{x}}(\mathbf{z}) + 2\beta (h(\mathbf{z}) - \nu) \pi_{\mathbf{x}}(\mathbf{z}) \stackrel{!}{=} 0 \quad (62)$$

$$\Rightarrow h^*(\mathbf{z}) = -\frac{\pi_{\mathbf{x}'}(\mathbf{z})}{2\beta \pi_{\mathbf{x}}(\mathbf{z})} + \frac{\alpha}{2\beta} + \nu. \quad (63)$$

Substituting into Eq. 59 and simplifying leaves us with the dual problem

$$\max_{\alpha, \beta \geq 0} \alpha \mu - \beta \zeta - \frac{\alpha^2}{4\beta} + \frac{\alpha}{2\beta} - \alpha \nu + \nu - \frac{1}{4\beta} \sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})}. \quad (64)$$

In the following, let us use $\rho = \sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})}$ as a shorthand for the expected likelihood ratio. The problem is concave in α . We can thus find the optimum α^* by setting the derivative to zero, which gives us $\alpha^* = 2\beta(\mu - \nu) + 1$. Because $\beta \geq 0$ and our theorem assumes that $\nu \leq \mu$, the value α^* is a feasible solution to the dual problem. Substituting into Eq. 64 and simplifying results in

$$\max_{\beta \geq 0} \alpha^* \mu - \beta \zeta - \frac{\alpha^{*2}}{4\beta} + \frac{\alpha^*}{2\beta} - \alpha^* \nu + \nu - \frac{1}{4\beta} \rho \quad (65)$$

$$= \max_{\beta \geq 0} \beta ((\mu - \nu)^2 - \sigma^2) + \mu + \frac{1}{4\beta} (1 - \rho). \quad (66)$$

Lemma F.3 shows that the expected likelihood ratio ρ is always greater than or equal to 1. Lemma F.4 shows that $(\mu - \nu)^2 - \sigma^2 \leq 0$. Therefore Eq. 66 is concave in β . The optimal value of β can again be found by setting the derivative to zero:

$$\beta^* = \sqrt{\frac{1 - \rho}{4((\mu - \nu)^2 - \sigma^2)}}. \quad (67)$$

Recall that our theorem assumes $\sigma^2 \geq (\mu - \nu)^2$ and thus β^* is real valued. Substituting Eq. 67 into Eq. 66 shows that the maximum of our dual problem is

$$\mu + \sqrt{(1 - \rho)((\mu - \nu)^2 - \sigma^2)}. \quad (68)$$

By duality, this is a lower bound on our primal problem $\min_{h \in \mathbb{F}} \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x}')} [h(\mathbf{z})]$. We know that our prediction is certifiably robust, i.e. $f(\mathbf{x}) = y$, if $\min_{h \in \mathbb{F}} \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x}')} [h(\mathbf{z})] > 0.5$. So, in particular, our prediction is robust if

$$\mu + \sqrt{(1 - \rho)((\mu - \nu)^2 - \sigma^2)} > 0.5 \quad (69)$$

$$\Leftrightarrow \rho < 1 + \frac{1}{\sigma^2 - (\mu - \nu)^2} \left(\mu - \frac{1}{2} \right)^2 \quad (70)$$

$$\Leftrightarrow \sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})} < 1 + \frac{1}{\sigma^2 - (\mu - \nu)^2} \left(\mu - \frac{1}{2} \right)^2 \quad (71)$$

The last equivalence is the result of inserting the definition of the expected likelihood ratio ρ . \square

With Theorem 5.1 in place, we can certify robustness for arbitrary smoothing distributions, assuming we can compute the expected likelihood ratio. When we are working with discrete data and the smoothing distributions factorize, this can be done efficiently, as the two following base certificates for binary data demonstrate.

F.3.1 BERNOULLI SMOOTHING FOR PERTURBATIONS OF BINARY DATA

We begin by proving the base certificate presented in § 5. Recall that we use a smoothing distribution $\mathcal{F}(\mathbf{x}, \boldsymbol{\theta})$ with $\theta \in [0, 1]^{D_{\text{in}}}$ that independently flips the d 'th bit with probability θ_d , i.e. for $\mathbf{x}, \mathbf{z} \in \{0, 1\}^{D_{\text{in}}}$ and $\mathbf{z} \sim \mathcal{F}(\mathbf{x}, \boldsymbol{\theta})$ we have $\Pr[z_d \neq x_d] = \theta_d$.

Corollary F.5. Given an output $g_n : \{0, 1\}^{D_{\text{in}}} \rightarrow \Delta_{|\mathbb{Y}|}$ mapping to scores from the $(|\mathbb{Y}| - 1)$ -dimensional probability simplex, let $f_n(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \mathbb{E}_{\mathbf{z} \sim \mathcal{F}(\mathbf{x}, \boldsymbol{\theta})} [g_n(\mathbf{z})_y]$ be the corresponding smoothed classifier with $\boldsymbol{\theta} \in [0, 1]^{D_{\text{in}}}$. Given an input $\mathbf{x} \in \{0, 1\}^{D_{\text{in}}}$ and smoothed prediction $y_n = f_n(\mathbf{x})$, let $\mu = \mathbb{E}_{\mathbf{z} \sim \mathcal{F}(\mathbf{x}, \boldsymbol{\theta})} [g_n(\mathbf{z})_y]$ and $\zeta = \operatorname{Var}_{\mathbf{z} \sim \mathcal{F}(\mathbf{x}, \boldsymbol{\theta})} [g_n(\mathbf{z})_y]$. Then, $\forall \mathbf{x}' \in \mathbb{H}^{(n)} : f_n(\mathbf{x}') = y_n$ with $\mathbb{H}^{(n)}$ defined as in Eq. 2, $w_d = \ln \left(\frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d} \right)$, $\eta = \ln \left(1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \right)$ and $p = 0$.

Proof. Based on our definition of the base certificates interface (see Definition 4.1, we must show that $\forall \mathbf{x}' \in \mathbb{H} : f_n(\mathbf{x}') = y_n$ with

$$\mathbb{H} = \left\{ \mathbf{x}' \in \{0, 1\}^{D_{\text{in}}} \mid \sum_{d=1}^{D_{\text{in}}} \ln \left(\frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d} \right) \cdot |x'_d - x_d|^0 < \ln \left(1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \right) \right\}, \quad (72)$$

Because all bits are flipped independently, our probability mass function $\pi_{\mathbf{x}}(\mathbf{z}) = \Pr_{\tilde{\mathbf{z}} \sim \Psi(\mathbf{x})} [\tilde{\mathbf{z}} = \mathbf{z}]$ factorizes:

$$\pi_{\mathbf{x}}(\mathbf{z}) = \prod_{d=1}^{D_{\text{in}}} \pi_{x_d}(z_d) \quad (73)$$

with

$$\pi_{x_d}(z_d) = \begin{cases} \theta_d & \text{if } z_d \neq x_d \\ 1 - \theta_d & \text{else} \end{cases}. \quad (74)$$

Thus, our expected likelihood ratio can be written as

$$\sum_{\mathbf{z} \in \{0, 1\}^{D_{\text{in}}}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})} = \sum_{\mathbf{z} \in \{0, 1\}^{D_{\text{in}}}} \prod_{d=1}^{D_{\text{in}}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} = \prod_{d=1}^{D_{\text{in}}} \sum_{z_d \in \{0, 1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)}. \quad (75)$$

For each dimension d , we can distinguish two cases: If both the perturbed and unperturbed input are the same in dimension d , i.e. $x'_d = x_d$, then $\frac{\pi_{x'_d}(z_d)}{\pi_{x_d}(z_d)} = 1$ and thus

$$\sum_{z_d \in \{0, 1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} = \sum_{z_d \in \{0, 1\}} \pi_{x'_d}(z_d) = \theta_d + (1 - \theta_d) = 1. \quad (76)$$

If the perturbed and unperturbed input differ in dimension d , then

$$\sum_{z_d \in \{0, 1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} = \frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d}. \quad (77)$$

Therefore, the expected likelihood ratio is

$$\prod_{d=1}^{D_{\text{in}}} \sum_{z_d \in \{0, 1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} = \prod_{d=1}^{D_{\text{in}}} \left(\frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d} \right)^{|x'_d - x_d|}. \quad (78)$$

Due to Theorem 5.1 (and using $\nu = \mu$ when computing the variance), we know that our prediction is robust, i.e. $f_n(\mathbf{x}') = y_n$, if

$$\sum_{\mathbf{z} \in \{0, 1\}^{D_{\text{in}}}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})} < 1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \quad (79)$$

$$\iff \prod_{d=1}^{D_{\text{in}}} \left(\frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d} \right)^{|x'_d - x_d|} < 1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \quad (80)$$

$$\iff \sum_{d=1}^{D_{\text{in}}} \ln \left(\frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d} \right) |x'_d - x_d| < \ln \left(1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \right). \quad (81)$$

Because x_d and x'_d are binary, the last inequality is equivalent to

$$\sum_{d=1}^{D_{\text{in}}} \ln \left(\frac{(1-\theta_d)^2}{\theta_d} + \frac{(\theta_d)^2}{1-\theta_d} \right) |x'_d - x_d|^0 < \ln \left(1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \right). \quad (82)$$

□

F.3.2 SPARSITY-AWARE SMOOTHING FOR PERTURBATIONS OF BINARY DATA

Sparsity-aware randomized smoothing (Bojchevski et al., 2020) is an alternative smoothing approach for binary data. It uses different probabilities for randomly deleting ($1 \rightarrow 0$) and adding ($0 \rightarrow 1$) bits to preserve data sparsity. For a random variable \mathbf{z} distributed according to the sparsity-aware distribution $\mathcal{S}(\mathbf{x}, \boldsymbol{\theta}^+, \boldsymbol{\theta}^-)$ with $\mathbf{x} \in \{0, 1\}^{D_{\text{in}}}$ and addition and deletion probabilities $\boldsymbol{\theta}^+, \boldsymbol{\theta}^- \in [0, 1]^{D_{\text{in}}}$, we have:

$$\begin{aligned}\Pr[z_d = 0] &= (1 - \theta_d^+)^{1-x_d} \cdot (\theta_d^-)^{x_d}, \\ \Pr[z_d = 1] &= (\theta_d^+)^{1-x_d} \cdot (1 - \theta_d^-)^{x_d}.\end{aligned}$$

The Bernoulli smoothing distribution we discussed in the previous section is a special case of sparsity-aware smoothing with $\boldsymbol{\theta}^+ = \boldsymbol{\theta}^-$. The runtime of the robustness certificate derived by Bojchevski et al. (2020) increases exponentially with the number of unique values in $\boldsymbol{\theta}^+$ and $\boldsymbol{\theta}^-$, which makes it unsuitable for localized smoothing. Variance-constrained smoothing, on the other hand, allows us to efficiently compute a certificate in closed form.

Corollary F.6. *Given an output $g_n : \mathbb{R}^{D_{\text{in}}} \rightarrow \Delta_{|\mathbb{Y}|}$ mapping to scores from the $(|\mathbb{Y}| - 1)$ -dimensional probability simplex, let $f_n(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \mathbb{E}_{\mathbf{z} \sim \mathcal{S}(\mathbf{x}, \boldsymbol{\theta}^+, \boldsymbol{\theta}^-)} [g_n(\mathbf{z})_y]$ be the corresponding smoothed classifier with $\boldsymbol{\theta}^+, \boldsymbol{\theta}^- \in [0, 1]^{D_{\text{in}}}$. Given an input $\mathbf{x} \in \{0, 1\}^{D_{\text{in}}}$ and smoothed prediction $y_n = f_n(\mathbf{x})$, let $\mu = \mathbb{E}_{\mathbf{z} \sim \mathcal{S}(\mathbf{x}, \boldsymbol{\theta}^+, \boldsymbol{\theta}^-)} [g_n(\mathbf{z})_y]$ and $\zeta = \operatorname{Var}_{\mathbf{z} \sim \mathcal{S}(\mathbf{x}, \boldsymbol{\theta}^+, \boldsymbol{\theta}^-)} [g_n(\mathbf{z})_y]$. Then, $\forall \mathbf{x}' \in \mathbb{H} : f_n(\mathbf{x}') = y_n$ for*

$$\mathbb{H} = \left\{ \mathbf{x}' \in \{0, 1\}^{D_{\text{in}}} \mid \sum_{d=1}^{D_{\text{in}}} \gamma_d^+ \cdot \mathbb{I}[x_d = 0 \neq x'_d] + \gamma_d^- \cdot \mathbb{I}[x_d = 1 \neq x'_d] < \eta \right\}, \quad (83)$$

where $\gamma^+, \gamma^- \in \mathbb{R}^{D_{\text{in}}}$, $\gamma_d^+ = \ln \left(\frac{(\theta_d^-)^2}{1 - \theta_d^+} + \frac{(1 - \theta_d^-)^2}{\theta_d^+} \right)$, $\gamma_d^- = \ln \left(\frac{(1 - \theta_d^+)^2}{\theta_d^-} + \frac{(\theta_d^+)^2}{1 - \theta_d^-} \right)$ and $\eta = \ln \left(1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \right)$.

Proof. Just like with the Bernoulli distribution we discussed in the previous section, all bits are flipped independently, meaning our probability mass function $\pi_{\mathbf{x}}(\mathbf{z}) = \Pr_{\tilde{\mathbf{z}} \sim \Psi(\mathbf{x})} [\tilde{\mathbf{z}} = \mathbf{z}]$ factorizes:

$$\pi_{\mathbf{x}}(\mathbf{z}) = \prod_{d=1}^{D_{\text{in}}} \pi_{x_d}(z_d) \quad (84)$$

with

$$\pi_{x_d}(z_d) = \begin{cases} \theta_d & \text{if } z_d \neq x_d \\ 1 - \theta_d & \text{else} \end{cases}. \quad (85)$$

As before, our expected likelihood ratio can be written as

$$\sum_{\mathbf{z} \in \{0, 1\}^{D_{\text{in}}}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})} = \sum_{\mathbf{z} \in \{0, 1\}^{D_{\text{in}}}} \prod_{d=1}^{D_{\text{in}}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} = \prod_{d=1}^{D_{\text{in}}} \sum_{z_d \in \{0, 1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)}. \quad (86)$$

We can now distinguish three cases. If both the perturbed and unperturbed input are the same in dimension d , i.e. $x'_d = x_d$, then $\frac{\pi_{x'_d}(z_d)}{\pi_{x_d}(z_d)} = 1$ and thus

$$\sum_{z_d \in \{0, 1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} = \sum_{z_d \in \{0, 1\}} \pi_{x'_d}(z_d) = 1. \quad (87)$$

If $x'_d = 1$ and $x_d = 0$, i.e. a bit was added, then

$$\sum_{z_d \in \{0, 1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} = \sum_{z_d \in \{0, 1\}} \frac{\pi_1(z_d)^2}{\pi_0(z_d)} = \frac{\pi_1(0)^2}{\pi_0(0)} + \frac{\pi_1(1)^2}{\pi_0(1)} = \frac{(\theta_d^-)^2}{1 - \theta_d^+} + \frac{(1 - \theta_d^-)^2}{\theta_d^+} \quad (88)$$

If $x'_d = 0$ and $x_d = 1$, i.e. a bit was deleted, then

$$\sum_{z_d \in \{0,1\}} \frac{\pi_{x'_d}(z)^2}{\pi_{x_d}(z)} = \sum_{z_d \in \{0,1\}} \frac{\pi_0(z_d)^2}{\pi_1(z_d)} = \frac{\pi_0(0)^2}{\pi_1(0)} + \frac{\pi_0(1)^2}{\pi_1(1)} = \frac{(1 - \theta_d^+)^2}{\theta_d^-} + \frac{(\theta_d^+)^2}{1 - \theta_d^-}. \quad (89)$$

Therefore, the expected likelihood ratio is

$$\prod_{d=1}^{D_{\text{in}}} \sum_{z_d \in \{0,1\}} \frac{\pi_{x'_d}(z_d)^2}{\pi_{x_d}(z_d)} \quad (90)$$

$$= \prod_{d=1}^{D_{\text{in}}} \left(\frac{(\theta_d^-)^2}{1 - \theta_d^+} + \frac{(1 - \theta_d^-)^2}{\theta_d^+} \right)^{\mathbb{I}[x_d=0 \neq x'_d]} \left(\frac{(1 - \theta_d^+)^2}{\theta_d^-} + \frac{(\theta_d^+)^2}{1 - \theta_d^-} \right)^{\mathbb{I}[x_d=1 \neq x'_d]} \quad (91)$$

$$= \prod_{d=1}^{D_{\text{in}}} \exp(\gamma_d^+)^{\mathbb{I}[x_d=0 \neq x'_d]} \cdot \exp(\gamma_d^-)^{\mathbb{I}[x_d=1 \neq x'_d]}. \quad (92)$$

In the last equation, we have simply used the shorthands γ_d^+ and γ_d^- defined in Corollary F.6. Due to Theorem 5.1 (and using $\nu = \mu$ when computing the variance), we know that our prediction is robust, i.e. $f_n(\mathbf{x}') = y_n$, if

$$\sum_{z \in \{0,1\}^{D_{\text{in}}}} \frac{\pi_{\mathbf{x}'}(z)^2}{\pi_{\mathbf{x}}(z)} < 1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \quad (93)$$

$$\iff \prod_{d=1}^{D_{\text{in}}} \exp(\gamma_d^+)^{\mathbb{I}[x_d=0 \neq x'_d]} \cdot \exp(\gamma_d^-)^{\mathbb{I}[x_d=1 \neq x'_d]} < 1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \quad (94)$$

$$\iff \sum_{d=1}^{D_{\text{in}}} \gamma_d^+ \cdot \mathbb{I}[x_d = 0 \neq x'_d] \cdot \gamma_d^- \cdot \mathbb{I}[x_d = 1 \neq x'_d] < \ln \left(1 + \frac{1}{\zeta} \left(\mu - \frac{1}{2} \right)^2 \right). \quad (95)$$

□

Use for collective certification. It should be noted that this certificate does not comply with our interface for base certificates (see Definition 4.1), meaning we can not directly use it to certify robustness to norm-bound perturbations using our collective linear program from Theorem 4.2. We can however use it to certify collective robustness to the more refined threat model used in (Schuchardt et al., 2021): Let the set of admissible perturbed inputs be $\mathbb{B}_{\mathbf{x}} = \left\{ \mathbf{x}' \in \{0,1\}^{D_{\text{in}}} \mid \sum_{d=1}^{D_{\text{in}}} [x_d = 0 \neq x'_d] \leq \epsilon^+ \wedge \sum_{d=1}^{D_{\text{in}}} [x_d = 1 \neq x'_d] \leq \epsilon^- \right\}$ with $\epsilon^+, \epsilon^- \in \mathbb{N}_0$ specifying the number of bits the adversary is allowed to add or delete. We can now follow the procedure outlined in § 3.2 to combine the per-prediction base certificates into a collective certificate for our new collective perturbation model. As discussed in, we can bound the number of predictions that are robust to simultaneous attacks by minimizing the number of predictions that are certifiably robust according to their base certificates:

$$\min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(\mathbf{x}') = y_n] \geq \min_{\mathbf{x}' \in \mathbb{B}_{\mathbf{x}}} \sum_{n \in \mathbb{T}} \mathbb{I}[\mathbf{x}' \in \mathbb{H}^{(n)}]. \quad (96)$$

Inserting the linear inequalities characterizing our perturbation model and base certificates results in:

$$\min_{\mathbf{x}' \in \{0,1\}^{D_{\text{in}}}} \sum_{n \in \mathbb{T}} \mathbb{I} \left[\sum_{d=1}^{D_{\text{in}}} \gamma_d^+ \cdot \mathbb{I}[x_d = 0 \neq x'_d] + \gamma_d^- \cdot \mathbb{I}[x_d = 1 \neq x'_d] < \eta^{(n)} \right] \quad (97)$$

$$\text{s.t.} \quad \sum_{d=1}^{D_{\text{in}}} [x_d = 0 \neq x'_d] \leq \epsilon^+, \quad \sum_{d=1}^{D_{\text{in}}} [x_d = 1 \neq x'_d] \leq \epsilon^-. \quad (98)$$

Instead of optimizing over the perturbed input \mathbf{x}' , we can define two vectors $\mathbf{b}^+, \mathbf{b}^- \in \{0, 1\}^{\mathcal{D}_{\text{in}}}$ that indicate in which dimension bits were added or deleted. Using these new variables, Eq. 97 can be rewritten as

$$\min_{\mathbf{b}^+, \mathbf{b}^- \in \{0, 1\}^{\mathcal{D}_{\text{in}}}} \sum_{n \in \mathbb{T}} \mathbb{I} \left[(\gamma^+)^T \mathbf{b}^+ + (\gamma^-)^T \mathbf{b}^- < \eta^{(n)} \right] \quad (99)$$

$$\text{s.t. } \text{sum}\{\mathbf{b}^+\} \leq \epsilon^+, \quad \text{sum}\{\mathbf{b}^-\} \leq \epsilon^-, \quad (100)$$

$$\sum_{d|x_d=1} b_d^+ = 0, \quad \sum_{d|x_d=0} b_d^- = 0. \quad (101)$$

The last two constraints ensure that bits can only be deleted where $x_d = 1$ and bits can only be added where $x_d = 0$. Finally, we can use the procedure for replacing the indicator functions with indicator variables that we discussed in § D to restate the above problem as the mixed-integer problem

$$\min_{\mathbf{b}^+, \mathbf{b}^- \in \{0, 1\}^{\mathcal{D}_{\text{in}}}, \mathbf{t} \in \{0, 1\}^{\mathcal{D}_{\text{out}}}} \sum_{n \in \mathbb{T}} t_n \quad (102)$$

$$\text{s.t. } (\gamma^+)^T \mathbf{b}^+ + (\gamma^-)^T \mathbf{b}^- \geq (1 - t_n) \eta^{(n)}, \quad (103)$$

$$\text{sum}\{\mathbf{b}^+\} \leq \epsilon^+, \quad \text{sum}\{\mathbf{b}^-\} \leq \epsilon^-, \quad (104)$$

$$\sum_{d|x_d=1} b_d^+ = 0, \quad \sum_{d|x_d=0} b_d^- = 0. \quad (105)$$

The first constraint ensures that t_n can only be set to 0 if the l.h.s. is greater or equal η_n , i.e. only when the base certificate can no longer guarantee robustness. The efficiency of the certificate can be improved by applying any of the techniques discussed in § E.

F.3.3 GAUSSIAN SMOOTHING FOR PERTURBATIONS OF CONTINUOUS DATA

Even though we specifically proposed variance-constrained certification as a means of efficiently certifying anisotropically smoothed classifiers for discrete data, it can be generalized to continuous distributions by replacing sums with integrals and mass functions with density functions (the proof is analogous to that in Appendix F.3).

In the following, we assume Gaussian smoothing, i.e. $\Psi(\mathbf{x}) \sim \mathcal{N}(\mathbf{x}, \Sigma)$ with $\Sigma \in \mathbb{R}_+^{D \times D}$ with density function $\pi_{\mathbf{x}}$. In this case, the expected ratio between $\pi_{\mathbf{x}'}$ and $\pi_{\mathbf{x}}$ is the exponential of the squared Mahalanobis distance (see Table 2 of (Gil et al., 2013) with $\alpha = 2$), i.e.

$$\int_{\mathbb{R}^D} \frac{\pi_{\mathbf{x}'}(\mathbf{z})}{\pi_{\mathbf{x}}(\mathbf{z})} \pi_{\mathbf{x}'}(\mathbf{z}) d\mathbf{z} = \exp \left((\mathbf{x}' - \mathbf{x}) \Sigma^{-1} (\mathbf{x}' - \mathbf{x}) \right).$$

This leads us to the following corollary of Theorem 5.1:

Corollary F.7. *Given a function $h : \mathbb{R}^D \rightarrow \Delta_{|\mathbb{Y}|}$ mapping to scores from the $(|\mathbb{Y}| - 1)$ -dimensional probability simplex, let $f(\mathbf{x}) = \text{argmax}_{y \in \mathbb{Y}} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{x}, \Sigma)} [h(\mathbf{z})_y]$ with covariance matrix $\Sigma \in \mathbb{R}_+^{D \times D}$. Given an input $\mathbf{x} \in \mathbb{X}$ and smoothed prediction $y = f(\mathbf{x})$, let $\mu = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{x}, \Sigma)} [h(\mathbf{z})_y]$ and $\zeta = \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{x}, \Sigma)} \left[(h(\mathbf{z})_y - \nu)^2 \right]$ with $\nu \in \mathbb{R}$. Assuming $\nu \leq \mu$, then $f(\mathbf{x}') = y$ if*

$$(\mathbf{x}' - \mathbf{x}) \Sigma^{-1} (\mathbf{x}' - \mathbf{x}) < \ln \left(1 + \frac{1}{\zeta - (\mu - \nu)^2} \left(\mu - \frac{1}{2} \right) \right). \quad (106)$$

As with Theorem 5.1, The r.h.s. of Eq. 106 depends on the expected softmax score μ , a variable $\nu \leq \mu$ and the expected squared difference ζ between μ and ν . For $\nu = \mu$ the parameter ζ is the variance of the softmax score. A higher expected value and a lower variance allow us to certify robustness for larger adversarial perturbations.

For comparison, ANCER (Eiras et al., 2022) guarantees robustness for the smoothed prediction $y_n = \text{argmax}_{y \in \mathbb{Y}} \Pr [g(\mathbf{x}) = y]$ if

$$(\mathbf{x}' - \mathbf{x}) \Sigma^{-1} (\mathbf{x}' - \mathbf{x}) < \Phi^{-1}(q_{y_n})^2, \quad (107)$$

where q_{y_n} is the probability of predicting class y_n , i.e. $q_{y_n} = \Pr_{z \sim \mathcal{N}(\mathbf{x}, \Sigma)} [g(z) = y_n]$. Here, $g : \mathbb{R}^D \rightarrow \mathbb{Y}$ directly outputs a class label instead of a softmax score. We see that both the variance-constrained certificate and ANCER yield the same certified ellipsoid, scaled by a different factor. This factor is the certifiable radius η , i.e. the r.h.s. term of Eqs. (106) and (107). We also see that both certificates have the same computational complexity – they both involve calculation of the squared Mahalanobis distance and a constant number of operations for evaluation of the certifiable radius.

In the following, we briefly assess under which conditions which certificate yields a larger certifiable radius η . For this evaluation, we assume that $g(\mathbf{x}) = \operatorname{argmax}_y h(\mathbf{x})$, i.e. g predicts the class with the highest softmax score. We then vary the prediction probability q_{y_n} and the expected softmax score μ within $[0.5, 1.0]$. For each μ , we calculate the largest possible variance ζ (using the Bhatia–Davis inequality $\zeta \leq (1 - \mu) \cdot \mu$), which will give us the weakest possible variance-constrained certificate (see Eq. (106)).

Fig. 13 shows the difference in certifiable radius η , with the dashed line indicating parameters for which both certificates are identical. We have omitted all combinations of q_{y_n} and μ that are not possible, namely $\mu > q_{y_n} + \frac{1}{2}(1 - q_{y_n})$. We see that ANCER is stronger when q_{y_n} is large, i.e. almost all samples from the smoothing distribution are correctly classified, but not necessarily with high confidence. The variance-constrained certificate is stronger when q_{y_n} is smaller and μ is larger, i.e. some samples are misclassified but the correctly classified ones have high confidence. Note however, that this is the worst case for the variance-constrained certificate. For $\zeta \rightarrow 0$, much larger radii can be certified (see Fig. 14 and Eq. (106)).

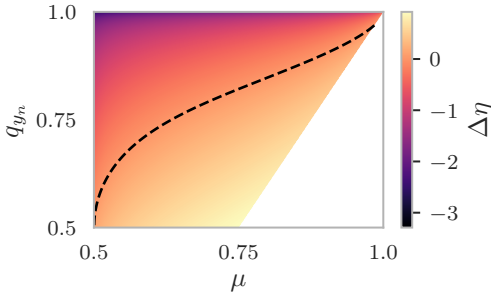


Figure 13: Worst-case difference in certifiable radius η between ANCER (Eiras et al., 2022) and the variance-constrained certificate for anisotropic Gaussian smoothing. The dashed line indicates combinations of prediction probability q_{y_n} and expected softmax score μ for which both certificates are equally strong.

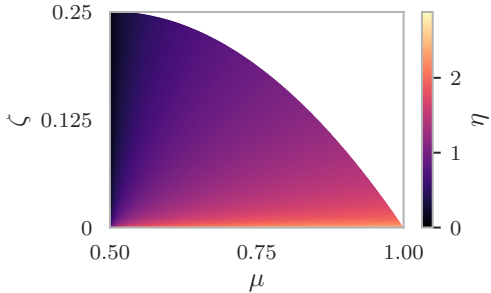


Figure 14: Certifiable radius η of the variance-constrained randomized smoothing certificate for anisotropic Gaussian smoothing as a function of the expected value μ and the variance ζ of the softmax score. If the variance is small, large radii can be certified – even if the expected softmax score is small.

G MONTE CARLO RANDOMIZED SMOOTHING

To make predictions and certify robustness, randomized smoothing requires computing certain properties of the distribution of a base model’s output, given an input smoothing distribution. For example, the certificate of Cohen et al. (2019) assumes that the smoothed model f predicts the most likely label output by base model g , given a smoothing distribution $\mathcal{N}(\mathbf{0}, \sigma \cdot \mathbf{1})$: $f(\mathbf{x}) = \operatorname{argmax}_{y \in \mathbb{Y}} \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma \cdot \mathbf{1})} [g(\mathbf{x} + \mathbf{z}) = y]$. To certify the robustness of a smoothed prediction $y = f(\mathbf{x})$ for a specific input \mathbf{x} , we have to compute the probability $q = \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma \cdot \mathbf{1})} [g(\mathbf{x} + \mathbf{z}) = y]$ to then calculate the maximum certifiable radius $\sigma \Phi^{-1}(q)$ with standard-normal inverse CDF Φ^{-1} . For complicated models like deep neural networks, computing such properties in closed form is usually not tractable. Instead, they have to be estimated using Monte Carlo sampling. The result are predictions and certificates that only hold with a certain probability.

Randomized smoothing with Monte Carlo sampling usually consists of three distinct steps:

1. First, a small number of samples N_1 from the smoothing distribution are used to generate a candidate prediction \hat{y} , e.g. the most frequently predicted class.
2. Then, a second round of N_2 samples is taken and a statistical test is used to determine whether the candidate prediction is likely to be the actual prediction of smoothed classifier f , i.e. whether $\hat{y} = f(\mathbf{x})$ with a certain probability $(1 - \alpha_1)$. If this is not the case, one has to abstain from making a prediction (or generate a new candidate prediction).
3. To certify the robustness of prediction \hat{y} , a final round of N_3 samples is taken to estimate all quantities needed for the certificate.

In the case of (Cohen et al., 2019), we need to estimate the probability $q = \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma \cdot \mathbf{1})} [g(\mathbf{x} + \mathbf{z}) = \hat{y}]$ to compute the certificate $\sigma \Phi^{-1}(q)$, whose strength is monotonically increasing in q . To ensure that the certificate holds with high probability $(1 - \alpha_2)$, we have to compute a probabilistic lower bound $\underline{q} \leq q$. Instead of performing two separate round of sampling, one can also re-use the same samples for the abstention test and certification. One particularly simple abstention mechanism is to just compute the Monte Carlo randomized smoothing certificate to determine whether $\forall \mathbf{x}' \in \{\mathbf{x}\} : f(\mathbf{x}') = \hat{y}$ with high probability, i.e. whether the prediction is robust to input \mathbf{x}' that is the result of ”perturbing” clean input \mathbf{x} with zero adversarial budget.

In the following, we discuss how we perform Monte Carlo randomized smoothing for our base certificates, as well as the baselines we use for our experimental evaluation. In § G.4, we discuss how we account for the multiple comparisons problem, i.e. the fact that we are not just trying to probabilistically certify a single prediction, but multiple predictions at once.

G.1 MONTE CARLO BASE CERTIFICATES FOR CONTINUOUS DATA

For our base certificates for continuous data, we follow the approach we already discussed in the previous paragraphs (recall that the certificate of Cohen et al. (2019) is a special case of our certificate with Gaussian noise for l_2 perturbations). We are given an input space $\mathbb{X}^{D_{\text{in}}}$, label space \mathbb{Y} , base model (or – in the case of multi-output classifiers – base model output) $g : \mathbb{X}^{D_{\text{in}}} \rightarrow \mathbb{Y}$ and smoothing distribution $\Psi(\mathbf{x})$ (either multivariate Gaussian or multivariate uniform). To generate a candidate prediction, we apply the base classifier to N_1 samples from the smoothing distribution in order to obtain predictions $(y^{(1)}, \dots, y^{(N_1)})$ and compute the majority prediction $\hat{y} = \operatorname{argmax}_{y \in \mathbb{Y}} \{n \mid y^{(n)} = \hat{y}\}$. Recall that for Gaussian and uniform noise, our certificate guarantees $\forall \mathbf{x}' \in \mathbb{H} : f(\mathbf{x}') = \hat{y}$ for

$$\mathbb{H} = \left\{ \mathbf{x}' \in \mathbb{X}^{D_{\text{in}}} \mid \sum_{d=1}^{D_{\text{in}}} w_d \cdot |x'_d - x_d|^p < \eta \right\},$$

with $\eta = (\Phi^{-1}(q))^2$ or $\eta = \Phi^{-1}(q)$ (depending on the distribution), $q = \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \sigma \cdot \mathbf{1})} [g(\mathbf{x} + \mathbf{z}) = \hat{y}]$ and standard-normal inverse CDF Φ^{-1} . To obtain a probabilistic certificate that holds with high probability $1 - \alpha$, we need a probabilistic lower bound on η . Both η are monotonically increasing in q , i.e. we can bound them by finding a lower bound \underline{q} on q . For this, we

take N_2 more samples from the smoothing distribution and compute a Clopper-Pearson lower confidence bound (Clopper & Pearson, 1934) on q . For abstentions, we use the aforementioned simple mechanism: We test whether $\mathbf{x} \in \mathbb{H}$. Given the definition of \mathbb{H} , this is equivalent to testing whether

$$\begin{aligned} 0 &< \Phi^{-1}(q) \\ \iff \Phi(0) &< \underline{q} \\ \iff 0.5 &< \underline{q}. \end{aligned}$$

If $\underline{q} \leq 0.5$, we abstain.

G.2 MONTE CARLO VARIANCE-CONSTRAINED CERTIFICATION

For variance-constrained certification, we smooth a model’s softmax scores. That is, we are given an input space $\mathbb{X}^{D_{\text{in}}}$, label space \mathbb{Y} , base model (or – in the case of multi-output classifiers – base model output) $g : \mathbb{X}^{D_{\text{in}}} \rightarrow \Delta_{|\mathbb{Y}|}$ with $(|\mathbb{Y}| - 1)$ -dimensional probability simplex $\Delta_{|\mathbb{Y}|}$ and smoothing distribution $\Psi(\mathbf{x})$ (Bernoulli or sparsity-aware noise, in the case of binary data). To generate a candidate prediction, we apply the base classifier to N_1 samples from the smoothing distribution in order to obtain vectors $(\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(N_1)})$ with $\mathbf{s} \in \Delta_{|\mathbb{Y}|}$, compute the average softmax scores $\bar{\mathbf{s}} = \frac{1}{N_1} \sum_{n=1}^{N_1} \mathbf{s}$ and select the label with the highest score $\hat{y} = \arg \max_y \bar{s}_y$.

Recall that our certificate guarantees robustness if the optimal value of the following optimization problem is greater than 0.5:

$$\min_{h: \mathbb{X} \rightarrow \mathbb{R}} \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x}')} [h(\mathbf{z})] \quad (108)$$

$$\text{s.t. } \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [h(\mathbf{z})] \geq \mu, \quad \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [(h(\mathbf{z}) - \nu)^2] \leq \zeta, \quad (109)$$

with $\mu = \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [g(\mathbf{z})_{\hat{y}}]$, $\zeta = \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [(g(\mathbf{z})_{\hat{y}} - \nu)^2]$ and a fixed scalar $\nu \in \mathbb{R}$. To obtain a probabilistic certificate, we have to compute a probabilistic lower bound on the optimal value of the optimization problem. Because it is a minimization problem, this can be achieved by loosening its constraints, i.e. computing a probabilistic lower bound $\underline{\mu}$ on μ and a probabilistic upper bound $\bar{\zeta}$ on ζ .

Like in CDF-smoothing (Kumar et al., 2020), we bound the parameters using CDF-based non-parametric confidence intervals. Let $F(s) = \Pr_{\mathbf{z} \sim \Psi(\mathbf{x})} [g(\mathbf{z})_{\hat{y}} \leq s]$ be the CDF of $g_{\hat{y}}(Z)$ with $Z \sim \Psi(\mathbf{x})$. Define M thresholds $0\tau_1 \leq \tau_2 \dots, \tau_{M-1} \leq \tau_M \leq 1$ with $\forall m : \tau_m \in [0, 1]$. We then take N_2 samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N_2)}$ from the smoothing distribution to compute the empirical CDF $\tilde{F}(s) = \sum_{n=1}^{N_2} \mathbb{I} [g(\mathbf{z}^{(n)})_{\hat{y}} \leq s]$. We can then use the Dvoretzky-Keifer-Wolfowitz inequality (Dvoretzky et al., 1956) to compute an upper bound \hat{F} and a lower bound \underline{F} on the CDF of $g_{\hat{y}}$:

$$\underline{F}(s) = \max(\tilde{F}(s) - v, 0) \leq F(s) \leq \min(\tilde{F}(s) + v, 1) = \bar{F}(s), \quad (110)$$

with $v = \sqrt{\frac{\ln 2/\alpha}{2 \cdot N_2}}$, which holds with high probability $(1 - \alpha)$. Using these bounds on the CDF, we can bound $\mu = \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [g(\mathbf{z})_{\hat{y}}]$ as follows (Anderson, 1969):

$$\mu \geq \tau_M - \tau_1 \bar{F}(\tau_1) + \sum_{m=1}^{M-1} (\tau_{m+1} - \tau_m) \bar{F}(\tau_m). \quad (111)$$

The parameter $\zeta = \mathbb{E}_{\mathbf{z} \sim \Psi(\mathbf{x})} [(g(\mathbf{z})_{\hat{y}} - \nu)^2]$ can be bounded in a similar fashion. Define $\xi_0, \dots, \xi_M \in \mathbb{R}_+$ with:

$$\begin{aligned} \xi_0 &= \max_{\kappa \in [0, \tau_1]} ((\kappa - \nu)^2) \\ \xi_M &= \max_{\kappa \in [\tau_M, 1]} ((\kappa - \nu)^2) \\ \xi_m &= \max_{\kappa \in [\tau_m, \tau_{m+1}]} ((\kappa - \nu)^2) \quad \forall m \in \{1, \dots, M-1\}, \end{aligned} \quad (112)$$

i.e. compute the maximum squared distance to ν within each bin $[\tau_m, \tau_{m+1}]$. Then:

$$\zeta \leq \xi_0 F(\tau_1) + \xi_M (1 - F(\tau_M)) + \sum_{m=1}^{M-1} \xi_m (F(\tau_{m+1}) - F(\tau_m)) \quad (113)$$

$$= \xi_M + \sum_{m=1}^{M-1} (\xi_{m-1} - \xi_m) F(\tau_m) \quad (114)$$

$$\leq \xi_M + \sum_{m=1}^{M-1} (\xi_{m-1} - \xi_m) (\text{sgn}(\xi_{m-1} - \xi_m) \bar{F}(\tau_m) + (1 - \text{sgn}(\xi_{m-1} - \xi_m)) \underline{F}(\tau_m)) \quad (115)$$

with probability $(1 - \alpha)$. In the first inequality, we bound the expected squared distance from ν by assuming that the probability mass in each bin $[\tau_m, \tau_{m+1}]$ is concentrated at the farthest point from ν . The equality is a result of reordering the telescope sum. In the second inequality, we upper-bound the CDF where it is multiplied with a non-negative value and lower-bound it where it is multiplied with a negative value.

With the probabilistic bounds $\underline{\mu}$ and $\bar{\zeta}$ we can now – in principle – evaluate our robustness certificate, i.e. check whether

$$\sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})} < 1 + \frac{1}{\bar{\zeta} - (\underline{\mu} - \nu)^2} \left(\underline{\mu} - \frac{1}{2} \right)^2. \quad (116)$$

where the π are the probability mass functions of smoothing distributions $\Psi(\mathbf{x})$ and $\Psi(\mathbf{x}')$. But one crucial detail of Theorem 5.1 underlying the certificate was that it only holds for $\nu \leq \underline{\mu}$. To use the method with Monte Carlo sampling, one has to ensure that $\nu \leq \underline{\mu}$ by first computing $\underline{\mu}$ and then choosing some smaller ν .

In our experiments, we use an alternative method that allows us to use arbitrary ν : From our proof of Theorem 5.1 we know that the dual problem of Eq. 108 is

$$\max_{\alpha, \beta \geq 0} \alpha \underline{\mu} - \beta \bar{\zeta} - \frac{\alpha^2}{4\beta} + \frac{\alpha}{2\beta} - \alpha \nu + \nu - \frac{1}{4\beta} \sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})}, \quad (117)$$

Instead of trying to find an optimal α (which causes problems in subsequent derivations if $\nu \not\leq \underline{\mu}$), we can simply choose $\alpha = 1$. By duality, the result is still a lower bound on the primal problem, i.e. the certificate remains valid. The dual problem becomes

$$\max_{\beta \geq 0} \underline{\mu} - \beta \bar{\zeta} + \frac{1}{4\beta} - \frac{1}{4\beta} \sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})}. \quad (118)$$

The problem is concave in β (because the expected likelihood ratio is ≥ 1). Finding the optimal β , comparing the result to 0.5 and solving for the expected likelihood ratio, shows that a prediction is robust if

$$\sum_{\mathbf{z} \in \mathbb{X}} \frac{\pi_{\mathbf{x}'}(\mathbf{z})^2}{\pi_{\mathbf{x}}(\mathbf{z})} < 1 + \frac{1}{\bar{\zeta}} \left(\underline{\mu} - \frac{1}{2} \right)^2. \quad (119)$$

For our abstention mechanism, like in the previous section, we compute the certificate \mathbb{H} and then test whether $\mathbf{x} \in \mathbb{H}$. In the case of Bernoulli smoothing and sparsity-aware smoothing), this corresponds to testing whether

$$1 < \ln \left(1 + \frac{1}{\bar{\zeta}} \left(\underline{\mu} - \frac{1}{2} \right) \right) \quad (120)$$

$$\iff \underline{\mu} > \frac{1}{2}. \quad (121)$$

G.3 MONTE CARLO CENTER SMOOTHING

While we can not use center smoothing as a base certificate, we benchmark our method against it during our experimental evaluation. The generation of candidate predictions, the abstention mechanism and the certificate are explained in (Kumar & Goldstein, 2021). The authors allow multiple options for generating candidate predictions. We use the "beta minimum enclosing ball" with $\beta = 2$ that is based on pair-wise distance calculations.

G.4 MULTIPLE COMPARISONS PROBLEM

The first step of our collective certificate is to compute one base certificate for each of the D_{out} predictions of the multi-output classifier. With Monte Carlo randomized smoothing, we want all of these probabilistic certificates to simultaneously hold with a high probability $(1 - \alpha)$. But as the number of certificates increases, so does the probability of at least one of them being invalid. To account for this *multiple comparisons problem*, we use Bonferroni (Bonferroni, 1936) correction, i.e. compute each Monte Carlo certificate such that it holds with probability $(1 - \frac{\alpha}{n})$.

For base certificates that only depend on $q_n = \Pr_{\mathbf{z} \sim \Psi^{(n)}} [g_n(\mathbf{z}) = \hat{y}_n]$, i.e. the probability of the base classifier predicting a particular label \hat{y}_n under the smoothing distribution, one can also use the strictly better Holm correction (Holm, 1979). This includes our Gaussian and uniform smoothing certificates for continuous data. Holm correction is a procedure than can be used to correct for the multiple comparisons problem when performing multiple arbitrary hypothesis tests. Given N hypotheses, their p -values are ordered in ascending order p_1, \dots, p_N . Starting at $i = 1$, the i 'th hypothesis is rejected if $p_i < \frac{\alpha}{N+1-i}$, until one reaches an i such that $p_i \geq \frac{\alpha}{N+1-i}$.

Fischer et al. (2021) proposed to use Holm correction as part of their procedure for certifying that all (non-abstaining) predictions of an image segmentation model are robust to adversarial perturbations. In the following, we first summarize their approach and then discuss how Holm correction can be used for certifying our notion of collective robustness, i.e. certifying the number of robust predictions. As in § G.1, the goal is to obtain a lower bound \underline{q}_n on $q_n = \Pr_{\mathbf{z} \sim \Psi^{(n)}} [g_n(\mathbf{z}) = \hat{y}_n]$ for each of the D_{out} classifier outputs. Assume we take N_2 samples $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(N_2)}$ from the smoothing distribution. Let $\nu_n = \sum_{i=1}^{N_2} \mathbb{I}[g_n(\mathbf{z}^{(i)}) = \hat{y}_n]$ and let $\pi : \{1, \dots, D_{\text{out}}\} \rightarrow \{1, \dots, D_{\text{out}}\}$ be a bijection that orders the ν_n in descending order, i.e. $\nu_{\pi(1)} \geq \nu_{\pi(2)} \dots \geq \nu_{\pi(D_{\text{out}})}$. Instead of using Clopper-Pearson confidence intervals to obtain tight lower bounds on the q_n , Fischer et al. (2021) define a threshold $\tau \in [0.5, 1)$ and use Binomial tests to determine for which n the bound $\tau \leq q_n$ holds with high-probability. Let $\text{BinP}(\nu_n, N_2, \leq, \tau)$ be the p-value of the one-sided binomial test, which is monotonically decreasing in ν_n . Following the Holm correction scheme, the authors test whether

$$\text{BinP}(\nu_{\pi(k)}, N_2, \leq, \tau) < \frac{\alpha}{D_{\text{out}} + 1 - k} \quad (122)$$

for $k = 1, \dots, D_{\text{out}}$ until reaching a k^* for which the null-hypothesis can no longer be rejected, i.e. the p-value is g.e.q. $\frac{\alpha}{D_{\text{out}} + 1 - k^*}$. They then know that with probability $1 - \alpha$, the bound $\tau \leq q_n$ holds for all $n \in \{\pi(k) \mid k \in \{1, \dots, k^*\}\}$. For these outputs, they use the lower bound τ to compute robustness certificates. They abstain with all other outputs.

This approach is sensible when one is concerned with the least robust prediction from a set of predictions. But our collective certificate benefits from having tight robustness guarantees for each of the individual predictions. Holm correction can be used with arbitrary hypothesis tests. For instance, we can use a different threshold τ_n per output g_n , i.e. test whether

$$\text{BinP}(\nu_{\pi(k)}, N_2, \leq, \tau_{\pi(k)}) < \frac{\alpha}{D_{\text{out}} + 1 - k} \quad (123)$$

for $k = 1, \dots, D_{\text{out}}$. In particular, we can use

$$\tau_n = \sup_t \text{ s.t. } \text{BinP}(\nu_n, N_2, \leq, t) < \frac{\alpha}{D_{\text{out}} + 1 - \pi^{-1}(n)}, \quad (124)$$

i.e. choose the largest threshold such that the null hypothesis can still be rejected. Eq. 124 is the lower Clopper-Pearson confidence bound with significance $\frac{\alpha}{D_{\text{out}} + 1 - \pi^{-1}(n)}$. This means that, instead of performing hypothesis tests, we can obtain probabilistic lower bounds $\underline{q}_n \leq q_n$ by computing Clopper-Pearson confidence bounds with significance parameters $\frac{\alpha}{D_{\text{out}}}, \dots, \frac{\alpha}{1}$. The \underline{q}_n can then be used to compute the base certificates. Due to the definition of the τ_n , all of the null hypotheses are rejected, i.e. we obtain valid probabilistic lower bounds on all q_n . We can thus use the abstention mechanism from § G.1, i.e. only abstain if $\underline{q}_n \leq 0.5$.

H COMPARISON TO THE COLLECTIVE CERTIFICATE OF FISCHER ET AL. (2021)

Our collective certificate based on localized smoothing is designed to bound the number of simultaneously robust predictions. Fischer et al. (2021) designed SegCertify to determine whether all predictions are simultaneously robust. As discussed in § 3.2, their work is based on the naïve collective certification approach applied to isotropic Gaussian smoothing: They first certify each output independently, then count the number of certifiably robust predictions for a specific adversarial budget and then test whether the number of certifiably robust predictions equals the overall number of predictions. To obtain better guarantees in practical scenarios, they further propose to

- use Holm correction to address the multiple comparisons problem (see § G.4),
- Abstain at a higher rate to avoid “bad componets”, i.e. predictions y_n that have a low consistency $q_n = \Pr_{\mathbf{z} \sim \mathcal{N}(\mathbf{x}, \sigma)} [g(\mathbf{z}) = y]$ and thus very small certifiable radii.

A more technical summary of their method can be found in § G.4.

In the following, we discuss why our certificate can always offer guarantees that are at least as strong as SegCertify, both for our notion of collective robustness (number of robust predictions) and their notion of collective robustness (robustness of all predictions). In short, isotropic smoothing is a special case of localized smoothing and Holm correction can also be used for our base certificates. Before proceedings, please read the discussion on Monte Carlo base certificates and Clopper-Pearson confidence intervals in § G.1 and the multiple comparisons problem in § G.4.

A direct consequence of the results in § G.4 is that using Clopper-Pearson confidence intervals and Holm correction will yield stronger per-prediction robustness guarantees and lower abstention rates than the method of Fischer et al. (2021). The Clopper-Pearson-based method only abstains if one cannot guarantee that $q_n > 0.5$ with high probability, while their method abstains if one cannot guarantee that $q_n \geq \tau$ with $\tau \geq 0.5$ (or specific other predictions abstain). For all non-abstaining predictions, the Clopper-Pearson-based certificate will be at least as strong as the one obtained using a single threshold τ , as it computes the tightest bound for which the null hypothesis can still be rejected (see Eq. 124).

Consequently, when certifying our notion of collective robustness, i.e. determining *the number* of robust predictions given adversarial budget ϵ , a naïve collective robustness certificate (i.e. counting the number of predictions whose robustness are guaranteed by the base certificates) based on Clopper-Pearson bounds will also be stronger than the method of Fischer et al. (2021). It should however be noted that their method could potentially be used with other methods of family-wise error rate correction, although they state that “these methods do not scale to realistic segmentation problems” and do not discuss any further details.

Conversely, when certifying their notion of collective robustness, i.e. determining whether *all* non-abstaining predictions are robust given adversarial budget ϵ , the certificate based on Clopper-Pearson confidence bounds is also at least as strong as that of Fischer et al. (2021). To certify their notion of robustness, they iterate over all predictions and determine whether all non-abstaining predictions are certifiably robust, given ϵ . Naturally, as the Clopper-Pearson-based certificates are stronger, any prediction that is robust according to (Fischer et al., 2021) is also robust according to the Clopper-Pearson-based certificates. The only difference is that, for $\tau > 0.5$, their method will have more abstaining predictions. But, due to the direct correspondence of Clopper-Pearson confidence bounds and Binomial tests, we can modify our abstention mechanism to obtain exactly the same set of abstaining predictions: We simply have to use $\underline{q}_n \leq \tau$ instead of $\underline{q}_n \leq 0.5$ as our criterion.

Finally, it should be noted that our proposed collective certificate based on linear programming is at least as strong as the naïve collective certificate (see Eq. 1.1 and Eq. 1.2 in § 3.2). Thus, letting the set of targeted predictions \mathbb{T} be the set of all non-abstaining predictions and checking whether the collective certificate guarantees robustness for all of \mathbb{T} will also result in a certificate that is at least as strong as that of Fischer et al. (2021) in their setting.

I COMPARISON TO THE COLLECTIVE CERTIFICATE OF SCHUCHARDT ET AL. (2021)

In the following, we first present the collective certificate for binary graph-structured data proposed by Schuchardt et al. (2021) (see § I.1. We then show that, when using sparsity-aware smoothing distributions (Bojchevski et al., 2020) – the family of smoothing distributions used both in our work and that of Schuchardt et al. (2021) – our certificate subsumes their certificate. That is, our collective robustness certificate based on localized randomized smoothing can provide the same robustness guarantees (see § I.2).

I.1 THE COLLECTIVE CERTIFICATE

Their certificate assumes the input space to be $\mathbb{G} = \{0, 1\}^{N \times D} \times \{0, 1\}^{N \times N}$ – the set of undirected attributed graphs with N nodes and D attributes per node. The model is assumed to be a multi-output classifier $f : \mathbb{G} \rightarrow \mathbb{Y}^N$ that assigns a label from label set \mathbb{Y} to each of the nodes. Given an input graph $\mathcal{G} = (\mathbf{X}, \mathbf{A})$ and a corresponding prediction $\mathbf{y} = f(\mathcal{G})$, they want to certify collective robustness to a set of perturbed graphs $\mathbb{B} \subseteq \mathbb{G}$. The perturbation model \mathbb{B} is characterized by four scalar parameters $r_{\mathbf{X}}^+, r_{\mathbf{X}}^-, r_{\mathbf{A}}^+, r_{\mathbf{A}}^- \in \mathbb{N}_0$, specifying the number of bits the adversary is allowed to add ($0 \rightarrow 1$) and delete ($1 \rightarrow 0$) in the attribute and adjacency matrix, respectively. It can also be extended to feature additional constraints (e.g. per-node budgets). We discuss how these can be integrated after showing our main result. A formal definition of the perturbation model can be found in Section B of (Schuchardt et al., 2021).

The goal of their work is to certify collective robustness for a set of targeted nodes $\mathbb{T} \subseteq \{1, \dots, N\}$, i.e. compute a lower bound on

$$\min_{G' \in \mathbb{B}} \sum_{n \in \mathbb{T}} \mathbb{I}[f_n(G') = y_n]. \quad (125)$$

Their approach to obtaining this lower-bound shares the same high-level idea as ours (see § 3.2): Combining per-prediction base certificates and leveraging some notion of locality. But while our method uses localized randomized smoothing, i.e. smoothing different outputs with different non-i.i.d. smoothing distributions to obtain base certificates that encode locality, their method uses a-priori knowledge about the strict locality of the classifier f . A model is strictly local if each of its outputs f_n only operates on a well-defined subset of the input data. To encode this strict locality, Schuchardt et al. (2021) associate each output f_n with an indicator vector $\psi^{(n)}$ and an indicator matrix $\Psi^{(n)}$ that fulfill

$$\begin{aligned} \sum_{m=1}^N \sum_{d=1}^D \psi_m^{(n)} \mathbb{I}[X_{m,d} \neq X'_{i,j}] + \sum_{i=1}^N \sum_{j=1}^N \Psi_m^{(n)} \mathbb{I}[A_{m,d} \neq A'_{i,j}] &= 0 \\ \implies f_n(\mathbf{X}, \mathbf{A}) &= f_n(\mathbf{X}', \mathbf{A}'). \end{aligned} \quad (126)$$

for any perturbed graph $\mathcal{G}' = (\mathbf{X}', \mathbf{A}')$. Eq. 126 expresses that the prediction of output f_n remains unchanged if all inputs in its receptive field remain unchanged. Conversely, it expresses that perturbations outside the receptive field can be ignored. Unlike in our work, Schuchardt et al. (2021) describe their base certificates as sets in adversarial budget space. That is, some certification procedure is applied to each output f_n to obtain a set

$$\mathbb{K}^{(n)} \subseteq [r_{\mathbf{X}}^+] \times [r_{\mathbf{X}}^-] \times [r_{\mathbf{A}}^+] \times [r_{\mathbf{A}}^-] \quad (127)$$

with $[k] = \{0, \dots, k\}$. If $[c_{\mathbf{X}}^+ \ c_{\mathbf{X}}^- \ c_{\mathbf{A}}^+ \ c_{\mathbf{A}}^-]^T \in \mathbb{K}^{(n)}$, then prediction y_n is robust to any perturbed input with exactly $c_{\mathbf{X}}^+$ attribute additions, $c_{\mathbf{X}}^-$ attribute deletions, $c_{\mathbf{A}}^+$ edge additions and $c_{\mathbf{A}}^-$ edge deletions. A more detailed explanation can be found in Section 3 of (Schuchardt et al., 2021). Note that the base certificates only depend on the number of perturbations, not their location in the input. Only by combining them using the receptive field indicators from Eq. 126 can one obtain a collective certificate that is better than the naïve collective certificate (i.e. counting how many predictions are certifiably robust to the collective threat model). The resulting collective certificate

is

$$\min_{\mathbf{b}^+, \mathbf{b}^-, \mathbf{B}^+, \mathbf{B}^-} \sum_{n \in \mathbb{T}} \mathbb{I} \left[\left[(\boldsymbol{\psi}^{(n)})^T \mathbf{b}_{\mathbf{X}}^+ \quad (\boldsymbol{\psi}^{(n)})^T \mathbf{b}_{\mathbf{X}}^- \quad \sum_{i,j} \Psi_{i,j}^{(n)} \mathbf{B}_{i,j}^+ \quad \sum_{i,j} \Psi_{i,j}^{(n)} \mathbf{B}_{i,j}^- \right]^T \in \mathbb{K}^{(n)} \right] \quad (128)$$

$$\text{s.t.} \quad \sum_{m=1}^N b_m^+ \leq r_{\mathbf{X}}^+, \quad \sum_{m=1}^N b_m^- \leq r_{\mathbf{X}}^-, \quad \sum_{i=1}^N \sum_{j=1}^N B_{i,j}^+ \leq r_{\mathbf{A}}^+, \quad \sum_{i=1}^N \sum_{j=1}^N B_{i,j}^- \leq r_{\mathbf{A}}^-, \quad (129)$$

$$\mathbf{b}^+, \mathbf{b}^- \in \mathbb{N}_0^N \quad \mathbf{B}^+, \mathbf{B}^- \in \mathbb{N}_0^{N \times N}. \quad (130)$$

The variables defined in Eq. 130 model how the adversary allocates their adversarial budget, i.e. how many attributes are perturbed per node and which edges are modified. Eq. 129 ensures that this allocation is compliant with the collective threat model. Finally, in Eq. 128 the indicator vector and matrix $\boldsymbol{\psi}^{(n)}$ and $\boldsymbol{\Psi}^{(n)}$ are used to mask out any allocated perturbation budget that falls outside the receptive field of f_n before evaluating its base certificate.

To solve the optimization problem, Schuchardt et al. (2021) replace each of the indicator functions with binary variables and include additional constraints to ensure that they have value 1 i.f.f. the indicator function would have value 1. To do so, they define one linear constraint per point separating the set of certifiable budgets $\mathbb{K}^{(n)}$ from its complement $\bar{\mathbb{K}}^{(n)}$ in adversarial budget space (the "Pareto front" discussed in Section 3 of (Schuchardt et al., 2021)).

From the above explanation, the main drawbacks of this collective certificate compared to our localized randomized smoothing approach and corresponding collective certificate should be clear. Firstly, if the classifier f is not strictly local, i.e. the receptive field indicators $\boldsymbol{\psi}$ and $\boldsymbol{\Psi}$ only have non-zero entries, then all base certificates are evaluated using the entire collective adversarial budget. It thus degenerates to the naïve collective certificate. Secondly, even if the model is strictly local, each of the outputs may assign varying levels of importance to different parts of its receptive field. Their method is incapable of capturing this additional soft locality. Finally, their means of evaluating the base certificates may involve evaluating a large number of linear constraints. Our method, on the other hand, only requires a single constraint per prediction. Our collective certificate can thus be more efficiently computed.

I.2 PROOF OF SUBSUMPTION

In the following, we show that any robustness certificate obtained by using the collective certificate of Schuchardt et al. (2021) with sparsity-aware randomized smoothing base certificates can also be obtained by using our proposed collective certificate with an appropriately parameterized localized smoothing distribution. The fundamental idea is that, for randomly smoothed models, completely randomizing all input dimensions outside the receptive field is equivalent to masking out any perturbations outside the receptive field.

First, we derive the certificate of Schuchardt et al. (2021) for predictions obtained via sparsity-aware smoothing. Schuchardt et al. (2021) require base certificates that guarantee robustness when $[c_{\mathbf{X}}^+, c_{\mathbf{X}}^-, c_{\mathbf{A}}^+, c_{\mathbf{A}}^-]^T \in \mathbb{K}^{(n)}$, where the c indicate the number of added and deleted attribute and adjacency bits. That is, the certificates must only depend on the number of perturbations, not on their location. To achieve this, all entries of the attribute matrix and all entries of the adjacency matrix, respectively, must share the same distribution. For the attribute matrix, they define scalar distribution parameters $p_{\mathbf{X}}^+, p_{\mathbf{X}}^- \in [0, 1]$. Given attribute matrix $\mathbf{X} \in \{0, 1\}^{N \times D}$, they then sample random attribute matrices $\mathbf{Z}_{\mathbf{X}}$ that are distributed according to sparsity-aware smoothing distribution $\mathcal{S}(\mathbf{X}, \mathbf{1} \cdot p_{\mathbf{X}}^+, \mathbf{1} \cdot p_{\mathbf{X}}^-)$ (see § F.3.2), i.e.

$$\begin{aligned} \Pr[(Z_{\mathbf{X}})_{m,d} = 0] &= (1 - p_{\mathbf{X}}^+)^{1 - X_{m,d}} \cdot (p_{\mathbf{X}}^-)^{X_{m,d}}, \\ \Pr[(Z_{\mathbf{X}})_{m,d} = 1] &= (p_{\mathbf{X}}^+)^{1 - X_{m,d}} \cdot (1 - p_{\mathbf{X}}^-)^{X_{m,d}}. \end{aligned}$$

Given input adjacency matrix \mathbf{A} , random adjacency matrices $\mathbf{Z}_{\mathbf{A}}$ are sampled from the distribution $\mathcal{S}(\mathbf{A}, \mathbf{1} \cdot p_{\mathbf{A}}^+, \mathbf{1} \cdot p_{\mathbf{A}}^-)$. Applying Corollary F.6 (to the flattened and concatenated attribute and adjacency matrices) shows that smoothed prediction $y_n = f_n(\mathbf{X}, \mathbf{A})$ is robust to the perturbed graph

$(\mathbf{X}', \mathbf{A}')$ if

$$\begin{aligned} & \sum_{m=1}^N \sum_{d=1}^D \gamma_{\mathbf{X}}^+ \cdot \mathbf{I}[X_{m,d} = 0 \neq X'_{m,d}] + \gamma_{\mathbf{X}}^- \cdot \mathbf{I}[X_{m,d} = 1 \neq X'_{m,d}] \\ & + \sum_{i=1}^N \sum_{j=1}^N \gamma_{\mathbf{A}}^+ \cdot \mathbf{I}[A_{i,j} = 0 \neq A'_{i,j}] + \gamma_{\mathbf{A}}^- \cdot \mathbf{I}[A_{i,j} = 1 \neq A'_{i,j}] \\ & < \eta^{(n)} \end{aligned} \quad (131)$$

$$\begin{aligned} \text{with } \gamma_{\mathbf{X}}^+ &= \ln \left(\frac{(p_{\mathbf{X}}^-)^2}{1-p_{\mathbf{X}}^+} + \frac{(1-p_{\mathbf{X}}^-)^2}{p_{\mathbf{X}}^+} \right), \quad \gamma_{\mathbf{X}}^- = \ln \left(\frac{(1-p_{\mathbf{X}}^+)^2}{p_{\mathbf{X}}^-} + \frac{(p_{\mathbf{X}}^+)^2}{1-p_{\mathbf{X}}^-} \right), \quad \gamma_{\mathbf{A}}^+ = \\ & \ln \left(\frac{(p_{\mathbf{A}}^-)^2}{1-p_{\mathbf{A}}^+} + \frac{(1-p_{\mathbf{A}}^-)^2}{p_{\mathbf{A}}^+} \right), \quad \gamma_{\mathbf{A}}^- = \ln \left(\frac{(1-p_{\mathbf{A}}^+)^2}{p_{\mathbf{A}}^-} + \frac{(p_{\mathbf{A}}^+)^2}{1-p_{\mathbf{A}}^-} \right) \text{ and } \eta^{(n)} = \ln \left(1 + \frac{1}{\sigma^{(n)2}} \left(\mu^{(n)} - \frac{1}{2} \right)^2 \right), \end{aligned}$$

where $\mu^{(n)}$ is the mean and $\sigma^{(n)}$ is the variance of the base classifier's output distribution, given the input smoothing distribution. Since the indicator functions for each perturbation type in Eq. 131 share the same weights, Eq. 131 can be rewritten as

$$\gamma_{\mathbf{X}}^+ c_{\mathbf{X}}^+ + \gamma_{\mathbf{X}}^- c_{\mathbf{X}}^- + \gamma_{\mathbf{A}}^+ c_{\mathbf{A}}^+ + \gamma_{\mathbf{A}}^- c_{\mathbf{A}}^- \leq \eta^{(n)}, \quad (132)$$

where $c_{\mathbf{X}}^+, c_{\mathbf{X}}^-, c_{\mathbf{A}}^+, c_{\mathbf{A}}^-$ are the overall number of added and deleted attribute and adjacency bits, respectively. Eq. 132 matches the notion of base certificates defined by Schuchardt et al. (2021), i.e. it corresponds to a set $\mathbb{K}^{(n)}$ in adversarial budget space for which we provably know that prediction y_n is certifiably robust if $[c_{\mathbf{X}}^+, c_{\mathbf{X}}^-, c_{\mathbf{A}}^+, c_{\mathbf{A}}^-]^T \in \mathbb{K}^{(n)}$. When we insert the base certificate Eq. 132 into objective function Eq. 128, the collective certificate of Schuchardt et al. (2021) becomes equivalent to

$$\begin{aligned} \min_{\mathbf{b}^+, \mathbf{b}^-, \mathbf{B}^+, \mathbf{B}^-} \sum_{n \in \mathbb{T}} \mathbf{I} \left[\gamma_{\mathbf{X}}^+ \left(\boldsymbol{\psi}^{(n)} \right)^T \mathbf{b}_{\mathbf{X}}^+ + \gamma_{\mathbf{X}}^- \left(\boldsymbol{\psi}^{(n)} \right)^T \mathbf{b}_{\mathbf{X}}^- \right. \\ \left. + \gamma_{\mathbf{A}}^+ \sum_{i,j} \Psi_{i,j}^{(n)} \mathbf{B}_{i,j}^+ + \sum_{i,j} \gamma_{\mathbf{A}}^- \Psi_{i,j}^{(n)} \mathbf{B}_{i,j}^- \leq \eta^{(n)} \right] \end{aligned} \quad (133)$$

$$\text{s.t. } \sum_{m=1}^N b_m^+ \leq r_{\mathbf{X}}^+, \quad \sum_{m=1}^N b_m^- \leq r_{\mathbf{X}}^-, \quad \sum_{i=1}^N \sum_{j=1}^N B_{i,j}^+ \leq r_{\mathbf{A}}^+, \quad \sum_{i=1}^N \sum_{j=1}^N B_{i,j}^- \leq r_{\mathbf{A}}^-, \quad (134)$$

$$\mathbf{b}^+, \mathbf{b}^- \in \mathbb{N}_0^N \quad \mathbf{B}^+, \mathbf{B}^- \in \mathbb{N}_0^{N \times N}. \quad (135)$$

Next, we show that obtaining base certificates through localized randomized smoothing with appropriately chosen parameters and using these base certificates within our proposed collective certificate (see Theorem 4.2) will result in the same optimization problem. Instead of using the same smoothing distribution for all outputs, we use different distribution parameters for each one. For the n 'th output, we sample random attributes matrices from distribution $\mathcal{S} \left(\mathbf{X}, \boldsymbol{\Theta}_{\mathbf{X}}^+(n), \boldsymbol{\Theta}_{\mathbf{X}}^-(n) \right)$ with $\boldsymbol{\Theta}_{\mathbf{X}}^+(n), \boldsymbol{\Theta}_{\mathbf{X}}^-(n) \in [0, 1]^{N \times D}$. Note that, in order to avoid having to index flattened vectors, we overload the definition of sparsity-aware smoothing to allow for matrix-valued parameters. For example, the value $\Theta_{\mathbf{X},n,d}^+(n)$ indicates the probability of flipping the value of input attribute $X_{n,d}$ from 0 to 1 and the value $\Theta_{\mathbf{X},n,d}^-(n)$ indicates the probability of flipping the value of input attribute $X_{n,d}$ from 1 to 0. We choose the following values for these parameters:

$$\Theta_{\mathbf{X},m,d}^+(n) = \psi_m^{(n)} \cdot p_{\mathbf{X}}^+ + \left(1 - \psi_m^{(n)} \right) \cdot 0.5, \quad (136)$$

$$\Theta_{\mathbf{X},m,d}^-(n) = \psi_m^{(n)} \cdot p_{\mathbf{X}}^- + \left(1 - \psi_m^{(n)} \right) \cdot 0.5, \quad (137)$$

where $\boldsymbol{\psi}^{(n)}$ is the receptive field indicator vector defined in Eq. 126 and $p_{\mathbf{X}}^+, p_{\mathbf{X}}^- \in [0, 1]$ are the same flip probabilities we used for the certificate of Schuchardt et al. (2021). Due to this parameterization, attribute bits inside the receptive field are randomized using the same distribution as in the certificate of Schuchardt et al. (2021), while attribute bits outside are set to either

0 or 1 with equal probability. Similarly, we sample random adjacency matrices from distribution $\mathcal{S}(\mathbf{A}, \Theta_{\mathbf{A}}^{+(n)}, \Theta_{\mathbf{A}}^{-(n)})$ with $\Theta_{\mathbf{A}}^{+(n)}, \Theta_{\mathbf{A}}^{-(n)} \in [0, 1]^{N \times D}$ and

$$\Theta_{\mathbf{A}i,j}^{+(n)} = \Psi_{i,j}^{(n)} \cdot p_{\mathbf{A}}^+ + (1 - \Psi_{i,j}^{(n)}) \cdot 0.5, \quad (138)$$

$$\Theta_{\mathbf{A}u,j}^{-(n)} = \Psi_{i,j}^{(n)} \cdot p_{\mathbf{A}}^- + (1 - \Psi_{i,j}^{(n)}) \cdot 0.5, \quad (139)$$

where $\Psi^{(n)}$ is the receptive field indicator matrix defined in Eq. 126. Note that, since we only alter the distribution of bits outside the receptive field, the smoothed prediction $y_n = f_n(\mathbf{X}, \mathbf{A})$ will be the same as the one obtained via the smoothing distribution used by Schuchardt et al. (2021). Applying Corollary F.6 (to the flattened and concatenated attribute and adjacency matrices) shows that smoothed prediction $y_n = f_n(\mathbf{X}, \mathbf{A})$ is robust to the perturbed graph $(\mathbf{X}', \mathbf{A}')$ if

$$\begin{aligned} & \sum_{m=1}^N \sum_{d=1}^D \tau_{\mathbf{X}m,d}^+ \cdot \mathbb{I}[X_{m,d} = 0 \neq X'_{m,d}] + \tau_{\mathbf{X}m,d}^- \cdot \mathbb{I}[X_{m,d} = 1 \neq X'_{m,d}] \\ & + \sum_{i=1}^N \sum_{j=1}^N \tau_{\mathbf{A}i,j}^+ \cdot \mathbb{I}[A_{i,j} = 0 \neq A'_{i,j}] + \tau_{\mathbf{A}i,j}^- \cdot \mathbb{I}[A_{i,j} = 1 \neq A'_{i,j}] \\ & < \eta^{(n)}. \end{aligned} \quad (140)$$

Because we only changed the distribution outside the receptive field, the scalar $\eta^{(n)}$, which depends on the output distribution's mean and variance μ and σ will be the same as the one obtained via the smoothing scheme used by Schuchardt et al. (2021) et al. Due to Corollary F.6 and the definition of our smoothing distribution parameters in Eqs. (136) to (139), the scalars $\tau_{\mathbf{X}m,d}^+, \tau_{\mathbf{X}m,d}^-, \tau_{\mathbf{A}i,j}^+, \tau_{\mathbf{A}i,j}^-$ have the following values:

$$\tau_{\mathbf{X}m,d}^+ = \psi_m^{(n)} \cdot \gamma_{\mathbf{X}}^+ + (1 - \psi_m^{(n)}) \cdot 2 \cdot \ln \left(\frac{(1 - 0.5)^2}{0.5} + \frac{0.5^2}{1 - 0.5} \right) \quad (141)$$

$$\tau_{\mathbf{X}m,d}^- = \psi_m^{(n)} \cdot \gamma_{\mathbf{X}}^- + (1 - \psi_m^{(n)}) \cdot 2 \cdot \ln \left(\frac{(1 - 0.5)^2}{0.5} + \frac{0.5^2}{1 - 0.5} \right) \quad (142)$$

$$\tau_{\mathbf{A}i,j}^- = \Psi_{i,j}^{(n)} \cdot \gamma_{\mathbf{A}}^+ + (1 - \Psi_{i,j}^{(n)}) \cdot 2 \cdot \ln \left(\frac{(1 - 0.5)^2}{0.5} + \frac{0.5^2}{1 - 0.5} \right) \quad (143)$$

$$\tau_{\mathbf{A}i,j}^+ = \Psi_{i,j}^{(n)} \cdot \gamma_{\mathbf{A}}^- + (1 - \Psi_{i,j}^{(n)}) \cdot 2 \cdot \ln \left(\frac{(1 - 0.5)^2}{0.5} + \frac{0.5^2}{1 - 0.5} \right), \quad (144)$$

where the γ are the same weights as those of the base certificate Eq. 131 of Schuchardt et al. (2021). Inserting the above values of τ into the base certificate Eq. 140 and using the fact that $\ln \left(\frac{(1-0.5)^2}{0.5} + \frac{0.5^2}{1-0.5} \right) = \ln(1) = 0$ results in

$$\begin{aligned} & \sum_{m=1}^N \sum_{d=1}^D \psi_m^{(n)} \cdot \gamma_{\mathbf{X}}^+ \cdot \mathbb{I}[X_{m,d} = 0 \neq X'_{m,d}] + \psi_m^{(n)} \cdot \gamma_{\mathbf{X}}^- \cdot \mathbb{I}[X_{m,d} = 1 \neq X'_{m,d}] \\ & + \sum_{i=1}^N \sum_{j=1}^N \Psi_{i,j}^{(n)} \cdot \gamma_{\mathbf{A}}^- \cdot \mathbb{I}[A_{i,j} = 0 \neq A'_{i,j}] + \Psi_{i,j}^{(n)} \cdot \gamma_{\mathbf{A}}^+ \cdot \mathbb{I}[A_{i,j} = 1 \neq A'_{i,j}] \\ & < \eta^{(n)}. \end{aligned} \quad (145)$$

While our collective certificate derived in § 4 only considers one perturbation type, we have already discussed how to certify robustness to perturbation models with multiple perturbation types in § F.3.2: We use a different budget variable per input dimension and perturbation type. Furthermore, the attribute bits of each node share the same noise level. Therefore, we can use the method discussed in § E.3, i.e. use a single budget variable per node instead of using one per node and attribute.

Modelling our collective problem in this way, using Eq. 145 as our base certificates and rewriting the first two sums using inner products results in the optimization problem

$$\min_{\mathbf{b}^+, \mathbf{b}^-, \mathbf{B}^+, \mathbf{B}^-} \sum_{n \in \mathbb{T}} \mathbb{I} \left[\gamma_{\mathbf{X}}^+ \left(\boldsymbol{\psi}^{(n)} \right)^T \mathbf{b}_{\mathbf{X}}^+ + \gamma_{\mathbf{X}}^- \left(\boldsymbol{\psi}^{(n)} \right)^T \mathbf{b}_{\mathbf{X}}^- \right. \\ \left. + \gamma_{\mathbf{A}}^+ \sum_{i,j} \Psi_{i,j}^{(n)} \mathbf{B}_{i,j}^+ + \sum_{i,j} \gamma_{\mathbf{A}}^- \Psi_{i,j}^{(n)} \mathbf{B}_{i,j}^- \leq \eta^{(n)} \right] \quad (146)$$

$$\text{s.t.} \quad \sum_{m=1}^N b_m^+ \leq r_{\mathbf{X}}^+, \quad \sum_{m=1}^N b_m^- \leq r_{\mathbf{X}}^-, \quad \sum_{i=1}^N \sum_{j=1}^N B_{i,j}^+ \leq r_{\mathbf{A}}^+, \quad \sum_{i=1}^N \sum_{j=1}^N B_{i,j}^- \leq r_{\mathbf{A}}^-, \quad (147)$$

$$\mathbf{b}^+, \mathbf{b}^- \in \mathbb{N}_0^N \quad \mathbf{B}^+, \mathbf{B}^- \in \mathbb{N}_0^{N \times N}. \quad (148)$$

This optimization problem is identical to that of Schuchardt et al. (2021) from Eqs. (133) to (135). The only difference is in how these problems would be mapped to a mixed-integer linear program. We would directly model the indicator functions in the objective using a single linear constraint. Schuchardt et al. (2021) would use multiple linear constraints, each corresponding to one point in the adversarial budget space.

To summarize: For randomly smoothed models, masking out perturbations using a-priori knowledge about a model’s strict locality is equivalent to completely randomizing (here: flipping bits with probability 50%) parts of the input. While Schuchardt et al. (2021) only derived their certificate for binary data, it can also be applied to strictly local models for continuous data. Considering our certificates for Gaussian (Proposition F.1) and uniform (Proposition F.2) smoothing, where the base certificate weights are $\frac{1}{\sigma^2}$ and $\frac{1}{\lambda}$, respectively, it is possible to perform the same masking operation as Schuchardt et al. (2021) by using $\sigma \rightarrow \infty$ and $\lambda \rightarrow \infty$.

Finally, it should be noted that the certificate by Schuchardt et al. (2021) allows for additional constraints, e.g. on the adversarial budget per node or the number of nodes controlled by the adversary. As all of them can be modelled using linear constraints on the budget variables (see Section C of their paper), they can be just as easily integrated into our mixed-integer linear programming certificate.