
Certified Adversarial Defense Methods for Inverse Problems

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 In this paper, we propose a randomized smoothing approach that aims to enhance
2 the robustness of the linear inverse problems against adversarial attacks, in par-
3 ticular rigorously guaranteeing an upper bound on a suitably-defined notion of
4 sensitivity to perturbations. In addition, we propose two novel algorithms that in-
5 incorporate randomized smoothing into training, where one algorithm injects random
6 perturbations to the input data directly, and the other algorithm adds random pertur-
7 bations to the gradients during backpropagation. We conduct numerical evaluations
8 on two of the most prominent inverse problems — denoising and compressed
9 sensing — utilizing a variety of neural network estimators and datasets. In broad
10 scenarios, these results demonstrate a strong potential of randomized smoothing
11 for enhancing the robustness of linear inverse problems.

12 1 Introduction

13 Linear inverse problems are fundamental in machine learning, signal processing, and statistics,
14 with the objective of recovering an unknown vector from an (often underdetermined) set of linear
15 measurements. Over the years, the incorporation of low-dimensional structures such as sparsity has
16 dominated the extensive research [12, 11, 15]. More recently, advances in deep learning have led to
17 their widespread adoption in inverse problems, including for signal modeling, decoder design, and
18 measurement design [32, 36]. We will specifically be interested in the decoder design aspect.

19 While deep learning methods for decoder design have shown promising results, recent studies have
20 suggested that they can be vulnerable to adversarial attacks that preserve the input’s semantics but
21 significantly worsen the accuracy of the output [20, 34, 3]. The study in [18] further indicates
22 that deep learning for inverse problems can come at the cost of instabilities, and current training
23 approaches cannot guarantee stable methods. Similar sensitivities of deep neural networks have been
24 well-studied in classification problems, where initial defense approaches were soon broken [6, 29].
25 This led to the development of *certified* defenses such as Reluplex [23], AI2 [16], Provable Defenses
26 [39], and randomized smoothing [8] have emerged, providing rigorous guarantees of robustness
27 against norm-bounded attacks in classification. However, there is much less work on the adversarial
28 robustness of inverse problems, which has fundamental differences from classification and even
29 regression, notably including the fact that the output is high-dimensional.

30 In this work, we propose a novel method to certify the robustness of linear inverse problems against
31 adversarial attacks using a suitably-adapted form of randomized smoothing [8]. Our approach is able
32 to provide certified upper bounds on a suitably-defined notion of sensitivity to perturbations. To our
33 knowledge, this kind of result has not been explored previously in the context of inverse problems.

34 In addition, inspired by analogies with classification and regression, we incorporate the concept of
35 randomized smoothing into the estimator’s training. We propose two versatile training algorithms that
36 can be applied to a variety of network architectures at least up to a moderate size. Through diverse

37 numerical experiments, we demonstrate that smoothing can provide significant robustness certificates
 38 against adversarial attacks, as well as being practically beneficial in both denoising and compressed
 39 sensing tasks.

40 **1.1 Setup**

41 A linear inverse problem consists of recovering a target signal¹ $\mathbf{x} \in \mathbb{R}^n$ from linear measurements
 42 $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ are the measurement matrix and additive noise,
 43 respectively. Given \mathbf{A} and \mathbf{y} , we employ a trainable estimator $f_{\theta} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ to approximate \mathbf{x} as
 44 $f_{\theta}(\mathbf{y})$, and assess correctness using the normalized mean squared error (MSE) criterion:

$$\ell(f_{\theta}(\mathbf{y}), \mathbf{x}) = \frac{1}{n} \|f_{\theta}(\mathbf{y}) - \mathbf{x}\|^2. \quad (1)$$

45 Parts of our paper are applicable to an arbitrary estimator f_{θ} , but at times we will also specifically
 46 consider training its parameters θ given access to a data set \mathcal{D} consisting of (\mathbf{x}, \mathbf{y}) pairs, with each \mathbf{y}
 47 being a measurement of the form $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ as above.

48 We study the adversarial robustness of f_{θ} by assuming the existence of an unknown adversarial
 49 attack ξ on the target signal \mathbf{x} at test time, such that the perturbed measurement vector is given by
 50 $\mathbf{y} = \mathbf{A}(\mathbf{x} + \xi) + \mathbf{b}$. To limit the power of the adversary, we adopt the standard assumption that ξ is
 51 ϵ -bounded in a ℓ_2 ball (i.e., $\|\xi\| \leq \epsilon$), and in the worst-case may be chosen to maximize the error
 52 $\ell(f_{\theta}(\mathbf{A}(\mathbf{x} + \xi) + \mathbf{b}), \mathbf{x})$. To streamline our notation, in the rest of this paper, we will represent the
 53 clean measurement $\mathbf{A}\mathbf{x} + \mathbf{b}$ as $\mathbf{y}_{\mathbf{x}}$ and represent the perturbed measurement $\mathbf{A}(\mathbf{x} + \xi) + \mathbf{b}$ as $\mathbf{y}_{\mathbf{x}+\xi}$,
 54 which can be generalized to $\mathbf{y}_{\mathbf{x}+\xi+\delta}$ in the case of multiple perturbations ξ and δ .

55 We note that our focus is on linear measurement since they are by far the most commonly considered
 56 in the literature on inverse problems, and already capture many interesting sub-problems such
 57 as inpainting, super-resolution, denoising, and compressive sensing. Nevertheless, extensions to
 58 non-linear models may be an interesting direction for future work.

59 **2 Related Work**

60 **Adversarial training** is a commonly used training strategy (primarily used for classification problems
 61 [25, 30, 33, 39]) to improve the robustness of perturbations through training with adversarial examples
 62 [17]. It follows the framework of empirical risk minimization (ERM) but replaces the input data with
 63 some ℓ_2 -bounded adversarial perturbations. Several methods have been proposed in prior work [25] to
 64 generate these perturbations, with the multi-step iterative method Projected Gradient Descent (PGD)
 65 being particularly effective. [22] studied the tradeoff between robustness and standard accuracy using
 66 adversarial training in linear regression problems, whereas [18] suggested that adversarial training
 67 alone may be insufficient for addressing instabilities in linear inverse problems. Nonetheless, we
 68 include adversarial training equipped with the PGD method as one of the baselines in our experiments
 69 and provide the ERM formulation as

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} [\ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi^*}), \mathbf{x})], \quad \text{where } \xi^* = \operatorname{argmax}_{\xi: \|\xi\|_2 \leq \epsilon} \ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi}), \mathbf{x}). \quad (2)$$

70 **Jacobian regularization** is a technique to improve robustness against adversarial attacks by penalizing
 71 large derivatives with respect to input data, specifically the norm of the input-output Jacobian
 72 matrix. The idea was first proposed by [19] but only implemented as a layer-wise regularization
 73 due to computational constraints. Later it was implemented at full scale in [37] and was made more
 74 efficient in [38, 40]. [21, 2] showed that both the Frobenius and spectral norms of the Jacobian matrix
 75 are experimentally effective in improving the robustness of neural networks in regression problems.

76 In the context of linear inverse problems, we can motivate the use of the Jacobian by considering the
 77 Taylor expansion of $f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi})$ with respect to the adversarial perturbation ξ :

$$f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi}) = f_{\theta}(\mathbf{y}_{\mathbf{x}} + \mathbf{A}\xi) = f_{\theta}(\mathbf{y}_{\mathbf{x}}) + \mathbf{J}(\mathbf{y}_{\mathbf{x}})\mathbf{A}\xi + O(\|\xi\|^2), \quad (3)$$

78 where $\mathbf{J}(\mathbf{y}_{\mathbf{x}}) \in \mathbb{R}^{n \times m}$ is the input-output Jacobian matrix. When the ℓ_2 -radius of the adversarial
 79 perturbation ξ is small enough, the higher-order terms can be ignored, and the stability of the

¹For image data, we think of \mathbf{x} as being its vectorized version, with n being the total number of pixels.

80 estimation is mainly determined by $\mathbf{J}(\mathbf{y}_x)\mathbf{A}$. Therefore, we can use Jacobian regularization by
 81 penalizing the spectral norm $\max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \|\mathbf{J}(\mathbf{y})\mathbf{A}\|$, as proposed by [2], and control the strength of
 82 regularization using a hyperparameter β . The ERM formulation for training is provided as

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} [\ell(f_{\theta}(\mathbf{y}), \mathbf{x})] + \beta \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \|\mathbf{J}(\mathbf{y})\mathbf{A}\|. \quad (4)$$

83 **Randomized smoothing** is a technique that was first introduced for certifying a classifier against
 84 ℓ_2 -bounded noise, potentially on large neural networks and on large images [26, 8, 28]. It works
 85 by adding random noise to input data and then classifying the noisy input to obtain a probability
 86 distribution over possible labels. Several works have shown that randomized smoothing guarantees
 87 that no possible attacks under a particular constraint could succeed in classification problems. [5]
 88 proposed a certificate for discrete data and demonstrated its usefulness in graph neural networks. [27]
 89 presented certified robustness against Wasserstein adversarial attacks. [13] covered parameterized
 90 transformations and certified robustness in the parameter space. [14] developed a certification method
 91 for image and point cloud segmentation. Furthermore, [35] embedded randomized smoothing into the
 92 classifier’s adversarial training framework, which inspired us to propose Algorithm 1 in this paper.

93 Randomized smoothing has also been extended to certain tasks beyond classification. [7] used it
 94 to certify the regression of the object detector’s bounding-box coordinates. [4] applied it to certify
 95 the accuracy score in watermarking problems and proposed a training framework from which we
 96 are inspired to propose Algorithm 2. These two works are the closest to ours, but there are some
 97 significant differences. In particular, in our case, the measurements are mapped to a high-dimensional
 98 signal manifold, whereas the the prior works consider outputting a single real number or only a few
 99 real numbers. In our setting, a naive strategy of summarizing the output via its loss (a single real
 100 number) is infeasible due to the loss being unknown at the decoder, as we discuss below.

101 3 Proposed Methods

102 3.1 Randomized smoothing for inverse problems during testing

103 We first review the necessary background of randomized smoothing for estimation in a generic
 104 *scalar-valued output* setting. Given an ordinary real-valued estimator $F : \mathbb{R}^n \rightarrow \mathbb{R}$, randomized
 105 smoothing is a method that constructs a new, smoothed estimator G from the base estimator F at
 106 testing time. Under *mean smoothing*, the smoothed estimator G takes as input $\mathbf{x} \in \mathbb{R}^n$ and returns
 107 the expectation of F under isotropic Gaussian noise perturbation δ on \mathbf{x} , i.e.,

$$G(\mathbf{x}) = \mathbb{E}_{\delta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [F(\mathbf{x} + \delta)], \quad (5)$$

108 where the noise level σ of δ controls the tradeoff between robustness and accuracy.

109 While (5) is the go-to approach for classification, its direct use for estimating *continuous-valued*
 110 quantities may result in skewed results and location-dependent bounds, as detailed in [7]. To
 111 overcome this limitation, they proposed a more suitable “percentile smoothing” (see (6) below) for
 112 one-dimensional regression tasks. In high-dimensional inverse problems, although incorporating
 113 randomized smoothing may seem feasible by regressing over the scalar estimation loss $\ell(f_{\theta}(\mathbf{y}_x), \mathbf{x})$,
 114 we cannot access the original signal \mathbf{x} and only have the measurement \mathbf{y}_x at hand (or more generally
 115 $\mathbf{y}_{\mathbf{x}+\xi}$). Hence, instead of considering the loss to \mathbf{x} itself, we consider the loss between two different
 116 reconstructions, one of which has an additional perturbation. Formally, we define the following.

117 **Definition 3.1.** Given an estimator $f_{\theta} : \mathbb{R}^m \rightarrow \mathbb{R}^n$, an adversarially perturbed measurement $\mathbf{y}_{\mathbf{x}+\xi}$,
 118 and a Gaussian noise vector $\delta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, we use δ to perturb \mathbf{x} further and obtain the measurement
 119 $\mathbf{y}_{\mathbf{x}+\xi+\delta} = \mathbf{y}_{\mathbf{x}+\xi} + \mathbf{A}\delta$. We refer to $\ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi+\delta}), f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi}))$ as the *discrepancy*, and consider the
 120 percentile smoothing of the discrepancy as a function of \mathbf{x} and ξ :

$$D_p(\mathbf{x}, \xi) = \inf \{t \in \mathbb{R} \mid \mathbb{P}[\ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi+\delta}), f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi})) \leq t] \geq p\}, \quad (6)$$

121 where p is the quantile value.

122 The underlying idea behind the *discrepancy* is that a small $D_p(\mathbf{x}, \xi)$ indicates that the estimation
 123 retains consistency despite the presence of Gaussian noise δ . This consistency is advantageous
 124 for adversarial robustness, allowing the addition of Gaussian noise to drown out an adversarial
 125 perturbation without impacting the reconstruction too much.

126 For concreteness, we often focus on $p = 0.5$, i.e., the median. Although $D_p(\mathbf{x}, \boldsymbol{\xi})$ will generally
 127 not admit a closed form, we can still approximate it by applying a Monte Carlo approach with a
 128 suitably-chosen confidence level (see Algorithm 3 in Appendix C).

129 The preceding ideas give rise to two estimators of \mathbf{x} that we will consider throughout the paper:

- 130 • For certified estimation, the Monte Carlo procedure with S samples gives an appropriate
 131 quantile value $q \in (0, 1)$; we compute the relevant S values of $\ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\boldsymbol{\xi}+\boldsymbol{\delta}}), f_{\theta}(\mathbf{y}_{\mathbf{x}+\boldsymbol{\xi}}))$ with
 132 randomly sampled $\boldsymbol{\delta}$, sort them, and take the one at position $q \cdot S$ (with integer rounding).
- 133 • We also consider median-based estimation, with the same procedure but simply $q = \frac{1}{2}$.

134 We emphasize that the Gaussian noise $\boldsymbol{\delta}$ for smoothing is intentionally introduced by the algorithm,
 135 whereas $\boldsymbol{\xi}$ is considered to be unknown and adversarial. We henceforth refer to them as the smoothing
 136 noise $\boldsymbol{\delta}$ and the adversarial noise $\boldsymbol{\xi}$.

137 3.2 Discrepancy certification

138 We follow Lemma 1 from [7] and propose the following corollary to bound the worst-case smoothed
 139 discrepancy (see Appendix A for the proof):

140 **Corollary 3.2.** *The median smoothed discrepancy with adversarial noise can be bounded by the
 141 smoothed discrepancy in the absence of adversarial noise as*

$$D_{0.5}(\mathbf{x}, \boldsymbol{\xi}) \leq D_{\Phi(\frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0}) \quad \forall \boldsymbol{\xi} \text{ s.t. } \|\boldsymbol{\xi}\| \leq \epsilon, \quad (7)$$

142 where Φ is the cumulative distribution function of unit normal distribution, and $\mathbf{0}$ is a zero vector.

143 Equation (7) states that a worst-case upper bound on the median smoothed discrepancy $D_{0.5}(\mathbf{x}, \boldsymbol{\xi})$ is
 144 the $\Phi(\frac{\epsilon}{\sigma})$ -th percentile smoothed discrepancy $D_{\Phi(\frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0})$, for arbitrary $\boldsymbol{\xi}'$ within an ϵ -ball around
 145 zero. This corollary guarantees that the sensitivity of \mathbf{x} (as measured by the discrepancy) is minimally
 146 affected by arbitrary shifts up to radius ϵ . We focus on the median for simplicity, but the proof reveals
 147 a more general form for any specified quantile value in $(0, 1)$.

148 It is worth noting that \mathbf{x} and $\boldsymbol{\xi}$ in Corollary 3.2 represent generic vectors whose assignments may
 149 vary depending on how the corollary is applied. In particular, suppose that we are in a scenario
 150 where the true signal is \mathbf{x}^* , and the true adversarial noise is $\boldsymbol{\xi}^*$, so the measurement vector is $\mathbf{y}_{\mathbf{x}^*+\boldsymbol{\xi}^*}$.
 151 Then we may apply Corollary 3.2 with $\mathbf{x} = \mathbf{x}^* + \boldsymbol{\xi}^*$, meaning that we are giving a guarantee on the
 152 discrepancy for all signals sufficiently close to $\mathbf{x}^* + \boldsymbol{\xi}^*$ (rather than sufficiently close to \mathbf{x}^* itself; we
 153 cannot work with \mathbf{x}^* directly since $\mathbf{y}_{\mathbf{x}^*}$ is not known to the algorithm).

154 We also emphasize that our approach does not directly ensure “good” or “bad” estimation but rather
 155 provides a guarantee on the discrepancy. This is analogous to classification, where a certificate may
 156 show that the decision is unaffected within a certain radius, without claiming it to be the correct
 157 decision in the first place. While this is a standard limitation of smoothing methods, we still consider
 158 using discrepancy calculations, Peak Signal-to-Noise Ratio (PSNR) values (a standard measure of
 159 accuracy with the definition provided in Section 4), and Structural Similarity (SSIM) values (an
 160 indicator of perceived similarity between estimations and ground truths) for experimental evaluations.

161 3.3 Embedding randomized smoothing into training

162 The previous subsection covers the use of randomized smoothing during testing, and applies to
 163 arbitrary f_{θ} (e.g., pre-trained). In this section, we provide algorithms that use randomized smoothing
 164 during training, i.e., selecting the parameters θ of f_{θ} based on training data.

165 The first algorithm, termed “Smoothed Adversarial Training” and building on [35], is an approach that
 166 roughly follows adversarial training (as discussed in Section 2) but uses randomized smoothing when
 167 generating the ℓ_2 -bounded adversarial perturbation. Specifically, adversarial training is performed
 168 using the perturbation noise $\boldsymbol{\xi}^*$ obtained as the expected value under the smoothing noise $\boldsymbol{\delta}$. Formally,
 169 the ERM with smoothing is performed as

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} [\ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\boldsymbol{\xi}^*}), \mathbf{x})], \quad \text{where } \boldsymbol{\xi}^* = \underset{\boldsymbol{\xi}: \|\boldsymbol{\xi}\|_2 \leq \epsilon}{\operatorname{argmax}} \ell\left(\mathbb{E}_{\boldsymbol{\delta} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [f_{\theta}(\mathbf{y}_{\mathbf{x}+\boldsymbol{\xi}+\boldsymbol{\delta}})], \mathbf{x}\right). \quad (8)$$

170 The algorithm computes ξ^* using a first-order method similar to PGD, but with the difference
 171 that the perturbed measurement $\mathbf{y}_{\mathbf{x}+\xi}$ is further smoothed in the neighborhood around $\mathbf{y}_{\mathbf{x}+\xi+\delta}$
 172 before being used in the loss function. It is difficult to evaluate ξ^* exactly, so we use Monte Carlo
 173 approximations with S i.i.d. smoothing noise samples, $\delta_1, \dots, \delta_S \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, and compute the
 174 estimated $\hat{f}_\theta(\mathbf{y}_{\mathbf{x}+\xi+\delta})$ using the sample mean

$$\hat{f}_\theta(\mathbf{y}_{\mathbf{x}+\xi+\delta}) = \frac{1}{S} \sum_{k=1}^S f_\theta(\mathbf{y}_{\mathbf{x}+\xi+\delta_k}). \quad (9)$$

175 See Algorithm 1 for a complete description. Although adversarial training is an empirical defense tech-
 176 nique that generally lacks provable guarantees, we will observe promising results in our experiments.
 177 The same goes for our second algorithm introduced below.

Algorithm 1 Smoothed Adversarial Training

```

input Training set  $\mathcal{D}$ ; Learning rate  $\alpha$ ; Estimator  $f_\theta$  and measurement matrix  $\mathbf{A}$ ; Adversarial noise bound  $\epsilon$ , iteration count  $T_{itr}$ , and step size  $r$ ; Smoothing noise level  $\sigma$  and noise sample count  $S$ ;
for  $1 \leq i \leq |\mathcal{D}|$  do
     $\xi_0 = \mathbf{0}$ 
    for  $1 \leq j \leq T_{itr}$  do
         $\Delta = \nabla_{\mathbf{x}} \ell(\hat{f}_\theta(\mathbf{y}_{\mathbf{x}+\xi_{j-1}+\delta}), \mathbf{x})|_{\mathbf{x}=\mathbf{x}_i}$ 
        (where  $\hat{f}_\theta(\mathbf{y}_{\mathbf{x}+\xi_{j-1}+\delta})$  is from (9))
         $\xi'_j = \xi_{j-1} + r * \Delta$ 
         $\xi_j = \operatorname{argmin}_{\xi: \|\xi\| \leq \epsilon} \|\xi - \xi'_j\|$ 
    end for
     $\xi^* = \xi_{T_{itr}}$ 
     $\theta = \theta - \alpha * \nabla_\theta \ell(f_\theta(\mathbf{y}_{\mathbf{x}_i+\xi^*}), \mathbf{x}_i)$ 
end for
return  $f_\theta$ 

```

178

179 The second algorithm, termed ‘‘Smoothed Gradient Training’’ and building on [4], is built upon the
 180 regular ERM training framework without a regularization term. To embed randomized smoothing,
 181 we augment the clean measurements with smoothing noise δ and compute the estimator’s gradients
 182 using an expectation with respect to δ . Specifically, for a given training pair (\mathbf{x}, \mathbf{y}) , the gradient \mathbf{g}_θ
 183 with respect to the estimator’s parameter θ is computed as

$$\mathbf{g}_\theta = \mathbb{E}_{\delta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [\nabla_\theta \ell(f_\theta(\mathbf{y}_{\mathbf{x}+\delta}), \mathbf{x})]. \quad (10)$$

184 To approximate the gradients in the equation, we use the Monte Carlo method, averaging gradients
 185 across S draws of i.i.d. smoothing noise. However, directly adding a large amount of noise during
 186 training can make it unstable. To overcome this issue, we adopt the approach proposed by [4] and
 187 incrementally increase the noise levels within each epoch. This is achieved by introducing a step
 188 count T_{step} and injecting smoothing noise incrementally with standard deviation σ ranging from
 189 $\frac{1}{T_{step}} \sigma_{max}$ to σ_{max} , where σ_{max} is the maximum noise level. A detailed description of the algorithm
 190 can be found in Algorithm 2.

191 We make several observations both training algorithms described above. Firstly, the underlying
 192 idea behind our training methods shares a resemblance with their counterparts in classification:
 193 By incorporating randomized smoothing during test time, the neural network is guided to main-
 194 tain stable behavior, generating similar reconstructions even in the presence of input per-
 195 turbations. Secondly, for Smoothed Adversarial Training, the perturbed measurement is the one being
 196 smoothed, while for Smoothed Gradient Training, the gradients themselves are being smoothed.
 197 Thirdly, these algorithms are generally adaptable to any differentiable neural network estimator f_θ ,
 198 allowing for a diverse range of choices in the numerical experiments. Fourthly, a possible alter-
 199 native approach is to smooth the estimation loss instead. For example, in (8), ξ^* is obtained as
 200 $\operatorname{argmax}_{\xi: \|\xi\|_2 \leq \epsilon} \mathbb{E}_{\delta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [\ell(f_\theta(\mathbf{y}_{\mathbf{x}+\xi+\delta}), \mathbf{x})]$; and in (10), the gradient is computed not before

201 but after the expectation, i.e., $\nabla_{\theta} \mathbb{E}_{\delta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [\ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\delta}), \mathbf{x})]$. However, loss smoothing appears
 202 to be less effective, because it takes the expected value in a one-dimensional manifold space and as a
 203 result, diminishes the “diversity” of the randomness.

204 4 Experiments

205 We conduct experiments in two different tasks: denoising, where $\mathbf{A} = \mathbf{I}$ and \mathbf{b} is Gaussian noise
 206 with level σ_b ; and compressed sensing, where \mathbf{A} is a random Gaussian matrix with compression
 207 ratio m/n , and the additive noise \mathbf{b} is set to zero. We consider training and testing sets from both the
 208 same data domain and different data domains (i.e., transfer learning), as is common in this line of
 209 works. Specifically, for denoising, we use the network DPDNN [10] as the estimator trained on the
 210 (grayscale converted) DIV2K dataset [1] and tested on the DIV2K validation set and BSD68 dataset
 211 [31], while for compressed sensing, we use the estimator ISTA-Net⁺⁺ [41] trained on the Train400
 212 dataset [24] and tested on the Set11 [24] and a subset of the (grayscale converted) ImageNet [9]
 213 datasets. All timings are for a single Nvidia GeForce RTX 3090. The code for all these experiments
 214 has been uploaded as a supplementary file.

215 4.1 Experimental settings

216 To assess the impact of randomized smoothing during training, we compare three baselines with our
 217 proposed training algorithms:

- 218 1. **Ordinary Training (Ord):** the standard training algorithm that does not incorporate
 219 any perturbation defense mechanisms, with the ERM formulation simply presented as
 220 $\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} [\ell(f_{\theta}(\mathbf{y}), \mathbf{x})]$;
- 221 2. **Adversarial Training (Adv):** the adversarial training algorithm equipped with PGD method as
 222 introduced in (2);
- 223 3. **Jacobian Regularization (Jcb):** the Jacobian regularization training algorithm using the spec-
 224 tral norm for penalization as detailed in (4), where we determine the value of hyperparameter β
 225 using Algorithm 4 from [2];
- 226 4. **Smoothed Adversarial Training (Smt-Adv)** from Algorithm 1;
- 227 5. **Smoothed Gradient Training (Smt-Grad)** from Algorithm 2.

228 Training parameters such as the number of epochs, batch size, and learning rate, are shared across all
 229 algorithms, whereas other parameters such as smoothing noise levels and adversarial attack iteration
 230 count are exclusive to certain algorithms (see Appendix B for further details). To evaluate the
 231 effectiveness of randomized smoothing at test time, we explore different choices of adversarial attack
 232 radii and various smoothing noise levels. Specifically, at test time, we apply the PGD method to
 233 produce adversarial attacks using a step size and number of iterations matching those used in the
 234 Adv and Smt-Adv algorithms during training. These hyperparameters (specified in Appendix B)
 235 are selected to ensure that the adversarial noise reaches the boundary of its ℓ_2 -radius ball. For the
 236 denoising task, we consider ℓ_2 radii $\epsilon < 15$, while for the compressed sensing task, we consider
 237 $\epsilon < 7$, avoiding using extreme values and preventing the estimates from being excessively poor. For
 238 smoothing, we examine noise levels ranging from 0 to 45 and use Monte Carlo approximation with
 239 50 samples (see Appendix B for a discussion on other choices and the sensitivity to varying S) and
 240 confidence level $c = 0.99$. Additionally, since our training process involves an additive noise level of
 241 $\sigma_b \leq 15$ and a compression ratio of $m/n \leq 0.5$, we restrict our testing to additive noise levels and
 242 compression ratios within these ranges.

243 We follow a common practice of rescaling images to have pixel values in $[0, 1]$. We then compute and
 244 analyze both the discrepancy (introduced in Section 3.1), PSNR, and SSIM. The general definition
 245 of PSNR depends on the maximum possible pixel value, but our $[0, 1]$ -normalization simplifies this
 246 to $\text{PSNR} = 10 \log_{10} \frac{1}{\text{MSE}}$, where the MSE is as defined in (1). The definition of SSIM is more
 247 complicated to state but is well-known, so is omitted here.

248 **4.2 Certificate evaluation**

249 We first examine the certificates of all types of trained estimators under varying ℓ_2 radii of the
 250 adversarial noise. Afterwards, we will investigate how adjusting the smoothing noise used in training
 251 can impact the certificate.

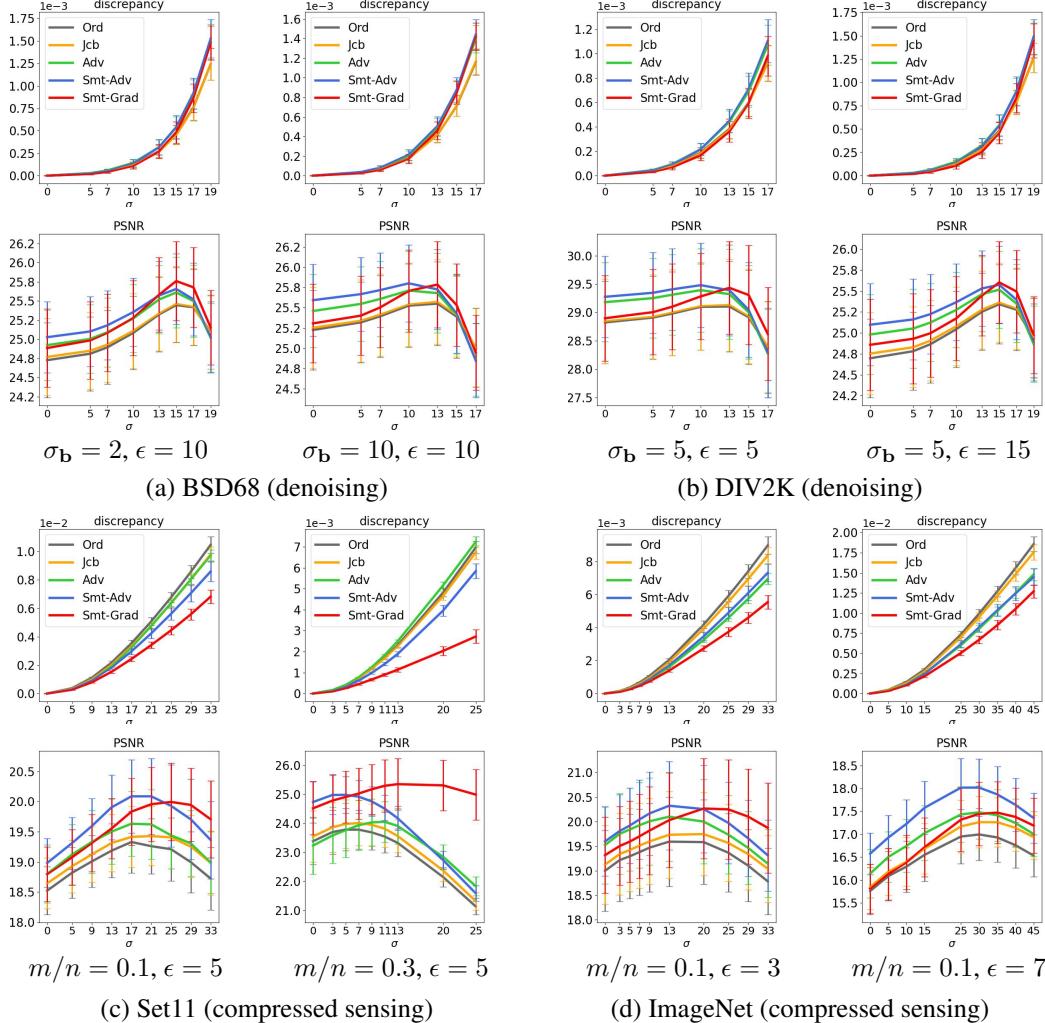


Figure 1: Certified smoothed discrepancy and its corresponding smoothed PSNR. The experimental settings for each pair of resulted discrepancy (top) and PSNR (bottom) are displayed in the footer.

252 Figure 1 presents a comparison of the certified results for both tasks. We display the results as the
 253 average over the corresponding dataset and include error bars representing half a standard deviation.
 254 We first observe that the certified upper bound of smoothed discrepancy increases with the smoothing
 255 level, which is expected as higher smoothing levels amount to more deviations from the original
 256 signal, as seen in (6). Regarding the PSNR of the reconstructions when smoothing is used, we
 257 consistently observe an ‘‘arc’’ shape in which the performance first improves with σ but then worsens,
 258 with the peak representing an ideal trade-off between σ being large enough to improve robustness, but
 259 not so large as to over-smooth. Naturally, the trend is that the best choice of σ increases with ϵ , with
 260 the two generally being comparable. However, ϵ is not the only factor at play; when the estimation
 261 problem is made ‘‘more difficult’’ via increased σ_b (denoising) or fewer measurements (compressed
 262 sensing), the best choice of σ also tends to increase.

263 In comparing the five training algorithms, we find that, at the optimal smoothing noise level, Smt-
 264 Grad consistently produces the lowest certified discrepancy and the highest corresponding estimation
 265 accuracy across all datasets and tasks. This suggests that Smt-Grad-trained estimators are the least

266 susceptible to perturbations, at least for the attacks we consider. Moreover, when little or no test-time
 267 smoothing is applied (small σ , seen in the left parts of the plots in Figure 1), Smt-Adv is generally
 268 seen to achieve the highest PSNR. On the other hand, when significant smoothing is applied, Smt-
 269 Grad can be preferable, particularly in the case of compressed sensing, with lower discrepancy and
 270 higher PSNR. Overall, we see that both of our proposed training algorithms can effectively improve
 271 the robustness of the estimator against adversarial attacks.

272 We now explore how different levels of smoothing noise σ used for training can affect performance
 273 in denoising and compressed sensing tasks. Table 1 provides a comparison of the performance
 274 of algorithms trained with different smoothing noise levels (shown as σ_{tr} to distinguish from the
 275 test-time smoothing noise σ) in terms of their performance metrics. These metrics include the
 276 PSNR of the clean estimate $f_{\theta}(y_x)$, the PSNR of the perturbed estimate $f_{\theta}(y_{x+\xi})$, and the PSNR
 277 of the smoothed estimate $f_{\theta}(y_{x+\xi+\delta})$ that provides the certified smoothed discrepancy. For each
 278 ℓ_2 -radius ϵ of the perturbation ξ , we use a zero test-time smoothing noise level to obtain the PSNR
 279 of the perturbed estimate and an empirically optimal test-time σ to obtain the smoothed PSNR.
 280 For comparison, we also include a baseline vanilla scenario that showcases the performance of the
 281 algorithm without any robustness considerations, i.e., Ord. The results demonstrate that a significant
 282 level of σ_{tr} (e.g., $\sigma_{\text{tr}} = 10$, which we use as the default training setting) can, in most cases, result in
 283 improved test-time robustness, indicated by higher PSNR values. However, as the estimator is trained
 284 more robustly against stronger adversarial attacks (when $\epsilon > 0$), the accuracy of the estimator’s clean
 285 estimation drops, as shown in the second column. An analogous trade-off between robustness and
 286 accuracy has also been observed in the adversarial robustness literature, e.g., [30], [4].

algorithm trained with σ_{tr}	$\epsilon=0$	ℓ_2 radius ϵ at test time					
		$\sigma = 0$	$\sigma = 10$	$\sigma = 0$	$\sigma = 14$	$\sigma = 0$	$\sigma = 18$
Ord	32.71	27.48	28.15	24.27	24.90	22.71	23.18
Smt-Adv, $\sigma_{\text{tr}} = 10$	31.87	27.85	28.44	24.56	25.13	22.97	23.36
Smt-Adv, $\sigma_{\text{tr}} = 5$	32.21	27.69	28.39	24.43	25.07	22.84	23.30
Smt-Adv, $\sigma_{\text{tr}} = 1$	32.40	27.48	28.27	24.26	24.97	22.71	23.23
Smt-Grad, $\sigma_{\text{tr}} = 10$	31.74	27.61	28.49	24.48	25.18	22.91	23.40
Smt-Grad, $\sigma_{\text{tr}} = 5$	31.87	27.61	28.39	24.43	25.11	22.86	23.36
Smt-Grad, $\sigma_{\text{tr}} = 1$	31.92	27.54	28.07	24.43	24.87	22.85	23.16

(a) $\sigma_b = 2$ on BSD68 (denoising)

algorithm trained with σ_{tr}	$\epsilon=0$	ℓ_2 radius ϵ at test time					
		$\sigma = 0$	$\sigma = 3$	$\sigma = 0$	$\sigma = 21$	$\sigma = 0$	$\sigma = 35$
Ord	27.79	24.64	24.72	18.52	19.28	15.56	16.79
Smt-Adv, $\sigma_{\text{tr}} = 10$	26.91	24.63	24.69	18.98	20.04	16.33	17.83
Smt-Adv, $\sigma_{\text{tr}} = 5$	27.16	24.41	24.52	18.62	19.70	15.67	17.11
Smt-Grad, $\sigma_{\text{tr}} = 10$	27.30	24.85	24.90	18.79	19.99	15.78	17.15
Smt-Grad, $\sigma_{\text{tr}} = 5$	27.48	24.77	24.83	18.65	19.31	15.56	16.51

(b) $m/n = 0.1$ on Set11 (compressed sensing)

Table 1: Trade-off between smoothing noise used during training and the resulting certificates during testing. The second column corresponds to the PSNR of the clean estimate, while the remaining cells indicate the PSNR of the perturbed estimate and the smoothed PSNR of the smoothed estimate.

287 4.3 Empirical smoothed median estimate

288 In practice, one may not specifically require certificates, and accordingly, it is of interest to study
 289 the performance when one instead adopts the simpler approach of using the median, as described in
 290 Section 3.1. To address this, we present Table 2 demonstrating the similarity between the experimental
 291 median and its corresponding certificate across various scenarios. These scenarios encompass different
 292 tasks and datasets, and we note that the rows of the table are not meant to be directly compared
 293 with one another. Overall, our results reveal that when the ideal test-time smoothing level is used,
 294 the evaluated discrepancy and PSNR values of the smoothed median estimate are not substantially
 295 different from those with certification.

296 We visually compare the smoothed estimate $f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi+\delta})$ corresponding to the smoothed median
 297 discrepancy with the perturbed estimate $f_{\theta}(\mathbf{y}_{\mathbf{x}+\xi})$ in Figure 2. In addition, We showcase the vanilla
 298 scenario where Ord is utilized and neither perturbation nor test-time smoothing is incorporated.
 299 Notably, we observe significant improvements in estimation accuracy and feature clarity in both
 300 denoising and compressed sensing tasks when smoothing is applied. Even in the absence of test-time
 301 smoothing ($\sigma = 0$), estimators trained with Smt-Adv or Smt-Grad consistently outperform the other
 302 three baselines. Further visual comparisons can be found in Figures 3 and 4 in Appendix D.

algorithm	task details (dataset)	smoothed median	smoothed certificate
Ord	$\sigma_b = 10, \epsilon = 10, \sigma = 13$ (BSD68)	25.523 (4.165 × 1e-4)	25.515 (4.174 × 1e-4)
Jcb	$m/n = 0.1, \epsilon = 5, \sigma = 17$ (Set11)	19.424 (3.295 × 1e-3)	19.410 (3.303 × 1e-3)
Adv	$\sigma_b = 5, \epsilon = 15, \sigma = 15$ (DIV2K)	25.514 (5.402 × 1e-4)	25.516 (5.410 × 1e-4)
Smt-Adv	$m/n = 0.1, \epsilon = 7, \sigma = 25$ (ImageNet)	18.116 (6.135 × 1e-3)	18.011 (6.138 × 1e-3)
smt-Grad	$\sigma_b = 2, \epsilon = 10, \sigma = 15$ (BSD68)	25.759 (4.812 × 1e-4)	25.716 (4.821 × 1e-4)

Table 2: Samples of the evaluated PSNR (and discrepancy in brackets) when using the smoothed median estimate vs. the smoothed estimate with certificate, for various tasks.

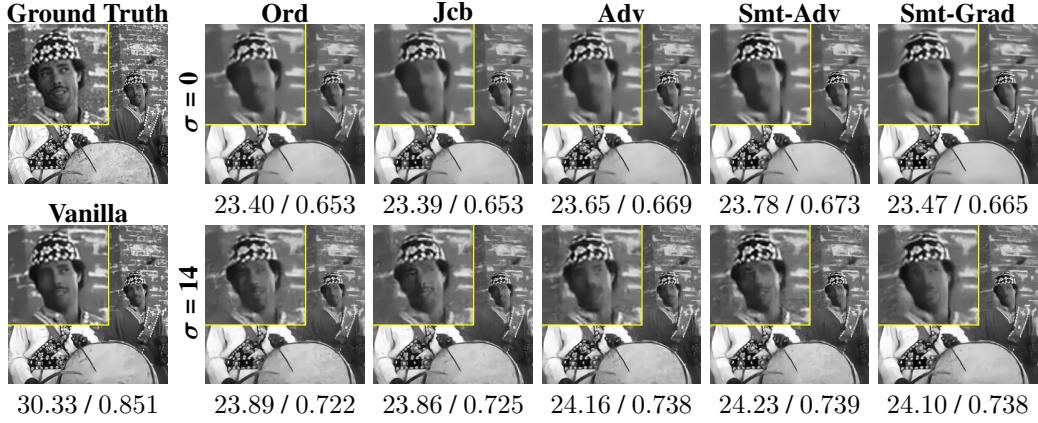
303 4.4 Note on computation

304 Our experiments have shown that the utilization of randomized smoothing during both training
 305 and testing phases of linear inverse problems effectively improves the robustness of the estimator
 306 against adversarial attacks. A slight caveat is that smoothing can increase the computation time, as
 307 is also the case with smoothing in other contexts. For training, we found that Smt-Grad was not
 308 too much slower than its unsmoothed counterpart, in part due to performing well with relatively few
 309 smoothing samples. On the other hand, the smoothing operation made Smt-Adv noticeably slower
 310 than Adv. Regarding test-time performance, determining the empirically ideal smoothing noise level,
 311 as depicted in Figure 1, can also be a time-consuming process.

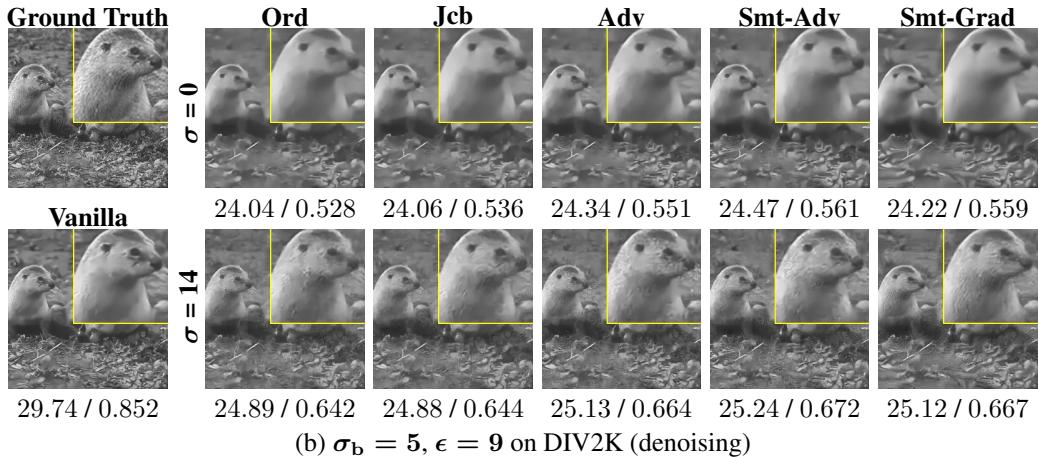
312 To support the above discussion, we present the training time per batch across all five algorithms in
 313 two tasks in Table 3 in Appendix B, where we observe that the training time required for Smt-Adv
 314 and Smt-Grad is higher than the baselines but still reasonable. We also provide Table 4 in Appendix B
 315 showcasing the testing time per data sample, where we observe that reducing smoothing samples from
 316 1000 to 50 leads to a significant speedup of around 15×, while maintaining comparable performance.

317 References

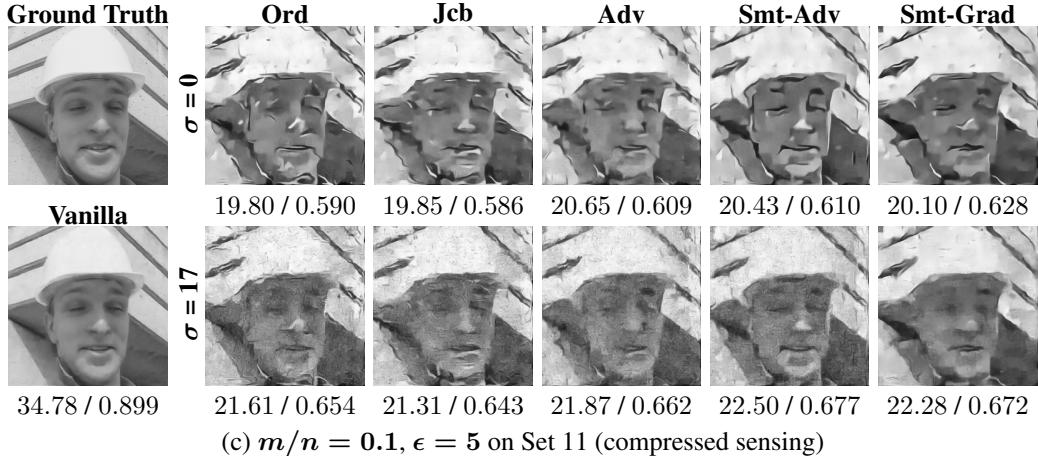
- 318 [1] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution:
 319 Dataset and study. In *IEEE Conference on Computer Vision and Pattern Recognition (Workshop)*,
 320 pages 1122–1131, 2017.
- 321 [2] Jaweria Amjad, Zhaoyan Lyu, and Miguel RD Rodrigues. Deep learning model-aware reg-
 322 ularization with applications to inverse problems. *IEEE Transactions on Signal Processing*,
 323 69:6371–6385, 2021.
- 324 [3] Vegard Antun, Francesco Renna, Clarice Poon, Ben Adcock, and Anders C Hansen. On
 325 instabilities of deep learning in image reconstruction and the potential costs of AI. *Proceedings
 326 of the National Academy of Sciences*, 117(48):30088–30095, 2020.
- 327 [4] Arpit Bansal, Ping-yeh Chiang, Michael J Curry, Rajiv Jain, Curtis Wigington, Varun Man-
 328 junatha, John P Dickerson, and Tom Goldstein. Certified neural network watermarks with
 329 randomized smoothing. In *International Conference on Machine Learning*, pages 1450–1465.
 330 PMLR, 2022.
- 331 [5] Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. Efficient robustness
 332 certificates for discrete data: Sparsity-aware randomized smoothing for graphs, images and
 333 more. In *International Conference on Machine Learning*, pages 1003–1013. PMLR, 2020.
- 334 [6] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In
 335 *IEEE Symposium on Security and Privacy*, pages 39–57, 2017.



(a) $\sigma_b = 2, \epsilon = 11$ on BSD68 (denoising)



(b) $\sigma_b = 5, \epsilon = 9$ on DIV2K (denoising)



(c) $m/n = 0.1, \epsilon = 5$ on Set 11 (compressed sensing)

Figure 2: Estimations (with the performance shown in the form of PSNR / SSIM) from the perturbed measurements in denoising and compressed sensing tasks. **Vanilla** represents the baseline unperturbed scenario where $\epsilon = 0, \sigma = 0$, and **Ord** is employed.

[7] Ping-yeh Chiang, Michael Curry, Ahmed Abdelkader, Aounon Kumar, John Dickerson, and Tom Goldstein. Detection as regression: Certified object detection with median smoothing. In *Advances in Neural Information Processing Systems*, volume 33, pages 1275–1286, 2020.

- 339 [8] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized
 340 smoothing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019.
- 341 [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
 342 hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*,
 343 pages 248–255, 2009.
- 344 [10] Weisheng Dong, Peiyao Wang, Wotao Yin, Guangming Shi, Fangfang Wu, and Xiaotong Lu.
 345 Denoising prior driven deep neural network for image restoration. *IEEE Transactions on Pattern
 346 Analysis and Machine Intelligence*, 41(10):2305–2318, 2018.
- 347 [11] Yonina C Eldar. *Sampling theory: Beyond bandlimited systems*. Cambridge University Press,
 348 2015.
- 349 [12] Yonina C Eldar and Gitta Kutyniok. *Compressed sensing: theory and applications*. Cambridge
 350 University Press, 2012.
- 351 [13] Marc Fischer, Maximilian Baader, and Martin Vechev. Certified defense to image transfor-
 352 mations via randomized smoothing. In *Advances in Neural Information Processing Systems*,
 353 volume 33, pages 8404–8417, 2020.
- 354 [14] Marc Fischer, Maximilian Baader, and Martin Vechev. Scalable certified segmentation via
 355 randomized smoothing. In *International Conference on Machine Learning*, pages 3340–3351.
 356 PMLR, 2021.
- 357 [15] Simon Foucart and Holger Rauhut. *A Mathematical Introduction to Compressive Sensing*.
 358 Birkhäuser Basel, 2013.
- 359 [16] Timon Gehr, Matthew Mirman, Dana Drachsler-Cohen, Petar Tsankov, Swarat Chaudhuri,
 360 and Martin Vechev. AI2: Safety and robustness certification of neural networks with abstract
 361 interpretation. In *IEEE Symposium on Security and Privacy*, pages 3–18, 2018.
- 362 [17] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversar-
 363 ial examples. *Stat*, 1050:20, 2015.
- 364 [18] Nina M. Gottschling, Vegard Antun, Anders C. Hansen, and Ben Adcock. The troublesome
 365 kernel – on hallucinations, no free lunches and the accuracy-stability trade-off in inverse
 366 problems. <https://arxiv.org/abs/2001.01258>, 2023.
- 367 [19] Shixiang Gu and Luca Rigazio. Towards deep neural network architectures robust to adversarial
 368 examples. In *International Conference on Learning Representations (Workshop)*, 2015.
- 369 [20] Yixing Huang, Tobias Würfl, Katharina Breininger, Ling Liu, Günter Lauritsch, and Andreas
 370 Maier. Some investigations on robustness of deep learning in limited angle tomography. In
 371 *Medical Image Computing and Computer Assisted Intervention*, pages 145–153, 2018.
- 372 [21] Daniel Jakubovitz and Raja Giryes. Improving dnn robustness to adversarial attacks using
 373 Jacobian regularization. In *European Conference on Computer Vision*, pages 514–529, 2018.
- 374 [22] Adel Javanmard, Mahdi Soltanolkotabi, and Hamed Hassani. Precise tradeoffs in adversarial
 375 training for linear regression. In *Conference on Learning Theory*, pages 2034–2078. PMLR,
 376 2020.
- 377 [23] Guy Katz, Clark Barrett, David L Dill, Kyle Julian, and Mykel J Kochenderfer. Reluplex:
 378 An efficient smt solver for verifying deep neural networks. In *International Conference on
 379 Computer Aided Verification*, pages 97–117, 2017.
- 380 [24] Kuldeep Kulkarni, Suhas Lohit, Pavan Turaga, Ronan Kerviche, and Amit Ashok. Reconnet:
 381 Non-iterative reconstruction of images from compressively sensed measurements. In *IEEE
 382 Conference on Computer Vision and Pattern Recognition*, pages 449–458, 2016.
- 383 [25] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale.
 384 In *International Conference on Learning Representations*, 2017.

- 385 [26] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified
 386 robustness to adversarial examples with differential privacy. In *IEEE Symposium on Security
 387 and Privacy*, pages 656–672, 2019.
- 388 [27] Alexander Levine and Soheil Feizi. Wasserstein smoothing: Certified robustness against
 389 wasserstein adversarial attacks. In *International Conference on Artificial Intelligence and
 390 Statistics*, pages 3938–3947. PMLR, 2020.
- 391 [28] Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness
 392 with additive noise. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 393 [29] Jiajun Lu, Hussein Sibai, Evan Fabry, and David Forsyth. No need to worry about adversarial
 394 examples in object detection in autonomous vehicles. In *IEEE Conference on Computer Vision
 395 and Pattern Recognition (Workshop)*, 2017.
- 396 [30] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
 397 Towards deep learning models resistant to adversarial attacks. In *International Conference on
 398 Learning Representations*, 2018.
- 399 [31] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human seg-
 400 mented natural images and its application to evaluating segmentation algorithms and measuring
 401 ecological statistics. In *IEEE International Conference on Computer Vision*, pages 416–423,
 402 2001.
- 403 [32] Gregory Ongie, Ajil Jalal, Christopher A Metzler, Richard G Baraniuk, Alexandros G Dimakis,
 404 and Rebecca Willett. Deep learning techniques for inverse problems in imaging. *IEEE Journal
 405 on Selected Areas in Information Theory*, 1(1):39–56, 2020.
- 406 [33] Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial
 407 examples. In *International Conference on Learning Representations*, 2018.
- 408 [34] Ankit Raj, Yoram Bresler, and Bo Li. Improving robustness of deep-learning-based image
 409 reconstruction. In *International Conference on Machine Learning*, pages 7932–7942. PMLR,
 410 2020.
- 411 [35] Hadi Salman, Jerry Li, Ilya Razenshteyn, Pengchuan Zhang, Huan Zhang, Sebastien Bubeck,
 412 and Greg Yang. Provably robust deep learning via adversarially trained smoothed classifiers. In
 413 *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 414 [36] Jonathan Scarlett, Reinhard Heckel, Miguel R. D. Rodrigues, Paul Hand, and Yonina C. Eldar.
 415 Theoretical perspectives on deep learning methods in inverse problems. *IEEE Journal on
 416 Selected Areas in Information Theory*, 3(3):433–453, 2022.
- 417 [37] Jure Sokolić, Raja Giryes, Guillermo Sapiro, and Miguel RD Rodrigues. Robust large margin
 418 deep neural networks. *IEEE Transactions on Signal Processing*, 65(16):4265–4280, 2017.
- 419 [38] Dániel Varga, Adrián Csizsárík, and Zsolt Zombori. Gradient regularization improves accuracy
 420 of discriminative models. *Schedae Informaticae*, 27, 2018.
- 421 [39] Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer
 422 adversarial polytope. In *International Conference on Machine Learning*, pages 5286–5295.
 423 PMLR, 2018.
- 424 [40] Yuichi Yoshida and Takeru Miyato. Spectral norm regularization for improving the generaliz-
 425 ability of deep learning. <https://arxiv.org/abs/1705.10941>, 2017.
- 426 [41] Di You, Jingfen Xie, and Jian Zhang. ISTA-Net++: flexible deep unfolding network for
 427 compressive sensing. In *IEEE International Conference on Multimedia and Expo*, pages 1–6,
 428 2021.

429 **Appendix**

430 **A Proof of Corollary 3.2**

431 **Corollary 3.2.** (Restated.) *The median smoothed discrepancy with adversarial noise can be bounded
432 by the smoothed discrepancy in the absence of adversarial noise as*

$$D_{0.5}(\mathbf{x}, \boldsymbol{\xi}) \leq D_{\Phi(\frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0}) \quad \forall \boldsymbol{\xi} \text{ s.t. } \|\boldsymbol{\xi}\| \leq \epsilon, \quad (11)$$

433 where Φ is the cumulative distribution function of unit normal distribution, and $\mathbf{0}$ is a zero vector.

434 *Proof.* According to Lemma 2.2 in [7], given a real-valued function $f : \mathbb{R}^m \rightarrow \mathbb{R}$, its percentile
435 smoothing function is defined as

$$\bar{h}_p(\mathbf{x}) = \inf\{t \in \mathbb{R} \mid \mathbb{P}[f(\mathbf{x} + \boldsymbol{\delta}) \leq t] \geq p\}, \quad (12)$$

436 and it is guaranteed that

$$\bar{h}_p(\mathbf{x} + \boldsymbol{\xi}) \leq \bar{h}_{\bar{p}}(\mathbf{x}) \quad \forall \|\boldsymbol{\xi}\|_2 \leq \epsilon, \quad (13)$$

437 where $\bar{p} = \Phi(\Phi^{-1}(p) + \frac{\epsilon}{\sigma})$.

438 Applying this result to our setting with any given estimator $f_{\theta} : \mathbb{R}^m \rightarrow \mathbb{R}^n$, given our definition of
439 the smoothed discrepancy

$$D_p(\mathbf{x}, \boldsymbol{\xi}) = \inf\{t \in \mathbb{R} \mid \mathbb{P}[\ell(f_{\theta}(\mathbf{y}_{\mathbf{x}+\boldsymbol{\xi}+\boldsymbol{\delta}}), f_{\theta}(\mathbf{y}_{\mathbf{x}+\boldsymbol{\xi}})) \leq t] \geq p\}, \quad (14)$$

440 we also have the guarantee

$$D_p(\mathbf{x}, \boldsymbol{\xi}) \leq D_{\bar{p}}(\mathbf{x}, \mathbf{0}) \quad \forall \boldsymbol{\xi} \text{ s.t. } \|\boldsymbol{\xi}\|_2 \leq \epsilon. \quad (15)$$

441 Then, by the definition of \bar{p} , we have

$$D_p(\mathbf{x}, \boldsymbol{\xi}) \leq D_{\Phi(\Phi^{-1}(p) + \frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0}) \quad \forall \boldsymbol{\xi} \text{ s.t. } \|\boldsymbol{\xi}\|_2 \leq \epsilon. \quad (16)$$

442 In particular, substituting 0.5 for p , we obtain

$$\begin{aligned} D_{0.5}(\mathbf{x}, \boldsymbol{\xi}) &\leq D_{\Phi(\Phi^{-1}(0.5) + \frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0}) \\ &= D_{\Phi(\frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0}) \end{aligned} \quad (17)$$

443 whenever $\|\boldsymbol{\xi}\|_2 \leq \epsilon$. □

444 **B Further Details on the Experiments**

445 We perform the experimental evaluation on both denoising and compressed sensing tasks. For each
446 task, we use a suitably-chosen neural network estimator equipped with one training set and two
447 testing sets. We consider training using all five algorithms discussed earlier. At test time, we apply the
448 PGD method to produce adversarial attacks using a step size and the number of iterations matching
449 those used in the Adv and Smt-Adv algorithms during training.

450 In addition, we select the training and testing parameters with minimal manual tuning, noting that
451 further fine-tuning may improve the performance. The details are given as follows:

- 452 • **DPDNN & DIV2K, BSD68** (denoising). We use the DPDNN model proposed in [10] as
453 the estimator f_{θ} . For training, we use the (grayscale converted) DIV2K dataset [1], which
454 consists of 1000 images of size 256×256 , and is further cropped to generate 32000 patches
455 of size 40×40 . For testing, we use the (grayscale converted) DIV2K validation set, which
456 contains 100 images cropped to a size of 384×384 , as well as the BSD68 dataset [31],
457 which consists of 68 images cropped to the 256×256 pixels in the center.

- 458 – Training parameters:

- 459 * (Shared across all algorithms) We follow the settings in [10] and set the number of
460 training epochs as 100, batch size as 32, learning rate $\alpha = 0.0005$, measurement
461 matrix $\mathbf{A} = \mathbf{I}$, and additive noise level $\sigma_b = 15$.

- 462 * (For Adv and Smt-Adv only) We set the iteration count $T_{itr} = 100$, step size
 463 $r = 0.4$, and adversarial noise ℓ_2 radius $\epsilon = 5$.
 464 * (For Smt-Adv and Smt-Grad only) We set the step count $T_{step} = 6$, sample count
 465 $S = 15$, and the smoothing noise level σ in Smt-Adv to be equal to the maximum
 466 smoothing noise level σ_{max} in Smt-Grad. Specifically, $\sigma = \sigma_{max} = 10$ (in Table 1
 467 we also use 1 and 5 for comparison).
 468 – Testing parameters:
 469 * (For measurement) $\mathbf{A} = \mathbf{I}$, and σ_b ranges from 0 to 15.
 470 * (For PGD attack) We set the attack iteration count $T_{itr} = 100$, step size $r = 0.4$,
 471 and consider multiple adversarial noise ℓ_2 radii with ϵ ranging from 0 to 15.
 472 * (For smoothing) We set the smoothing sample count $S = 50$ and consider multiple
 473 smoothing noise levels with σ ranging from 0 to 19.
 474 • **ISTA-Net⁺⁺ & Train400, Set11, ImageNet** (compressed sensing). The ISTA-Net⁺⁺
 475 model [41] is used as the estimator f_θ . We use the Train400 dataset [24], consisting of
 476 400 natural photos of size 180×180 for training, and the Set11 dataset [24] consisting
 477 of 11 grayscale images of size 256×256 for testing. Additionally, we use 300 randomly
 478 selected images sampled from the “tall building”, “open country”, and “street” categories
 479 from ImageNet [9] for testing. For each of these 300 images, the Y-channel luminance is
 480 extracted from the RGB color to form a grayscale-intensity image, and each is cropped to
 481 the central 256×256 pixels.
 482 – Training parameters:
 483 * (Shared across all algorithms) We follow the settings in [41] and set the number
 484 of training epoch as 100, batch size as 64, learning rate $\alpha = 0.0001$, additive
 485 noise level $\sigma_b = 0$, and Gaussian measurement matrices \mathbf{A} used in each batch
 486 with a randomly selected compression ratio from the set $\{0.1, 0.2, 0.3, 0.4, 0.5\}$.
 487 This randomization is done so that the neural network can learn to reconstruct at a
 488 variety of compression ratios.
 489 * (For Adv and Smt-Adv only) We set the iteration count $T_{itr} = 150$, step size
 490 $r = 0.2$, and adversarial noise ℓ_2 radius $\epsilon = 1$.
 491 * (For Smt-Adv and Smt-Grad only) We set the step count $T_{step} = 6$, step size
 492 $S = 15$, and the smoothing noise level σ in Smt-Adv to be equal to the maximum
 493 smoothing noise level σ_{max} in Smt-Grad. Specifically, $\sigma = \sigma_{max} = 10$ (in Table 1
 494 we also use 5 for comparison).
 495 – Testing parameters:
 496 * (For measurement) We use a Gaussian matrix \mathbf{A} with compression ratio ranging
 497 from 0.1 to 0.4 and no noise, i.e., $\mathbf{b} = \mathbf{0}$.
 498 * (For PGD attack) We set the iteration count $T_{itr} = 150$, step size $r = 0.2$, and
 499 consider multiple adversarial noise ℓ_2 radii with ϵ ranging from 0 to 7.
 500 * (For smoothing) We set the smoothing sample count $S = 50$ and consider multiple
 501 smoothing noise levels with σ ranging from 0 to 45.

502 Regarding the training time per batch across all five algorithms in two tasks, we present Table 3 below
 503 and observe that the training time required for Smt-Adv and Smt-Grad is higher than the baselines
 504 but still reasonable.

Estimator & Training Set	Ord	Jcb	Adv	Smt-Adv	Smt-Grad
DPDNN & DIV2K	0.49s	0.83s	3.53s	46.81s	45.48s
ISTA-Net ⁺⁺ & Train400	0.08s	0.15s	0.27s	2.49s	5.41s

Table 3: Training time per batch in denoising (top, batch size 32) and compressed sensing (bottom, batch size 64), where sample count of $S = 15$ for Smt-Adv and Smt-Grad is utilized in both tasks.

505 Regarding the choice of the number of Monte Carlo samples (S) during testing, Table 4 showcases
 506 the certified discrepancy and corresponding PSNR with S ranging from 50 to 1000 in chosen tasks.
 507 The results reveal that PSNR, in fact, shows very little degradation between smaller vs. higher values
 508 of S . For certified discrepancy, a downward trend is observed as expected (i.e., decreasing S gives
 509 a worse guarantee), yet not a particularly drastic one. Notably, lower S values significantly reduce
 510 the computation time as we see that reducing S from 1000 to 50 can give a roughly $15\times$ speedup.

511 Consequently, we opt for $S = 50$ in all experiments due to its consistent effectiveness, particularly
 512 with respect to PSNR (higher S may still be preferred if one specifically needs a strong certified
 513 guarantee).

	S=50 (t=2.20s)		S=500 (t=16.93s)		S=1000 (t=33.52s)	
	discrepancy	PSNR	discrepancy	PSNR	discrepancy	PSNR
Ord	5.085(± 0.948)	25.549(± 1.039)	5.083(± 0.945)	25.552(± 1.032)	5.072(± 0.923)	25.552(± 1.036)
Jcb	4.961(± 0.822)	25.574(± 1.045)	4.958(± 0.812)	25.572(± 1.045)	4.955(± 0.812)	25.573(± 1.044)
Adv	4.626(± 0.770)	25.671(± 0.980)	4.611(± 0.766)	25.679(± 0.971)	4.613(± 0.766)	25.678(± 0.972)
Smt-Adv	4.191(± 0.585)	25.724(± 0.987)	4.179(± 0.581)	25.727(± 0.981)	4.180(± 0.567)	25.727(± 0.982)
Smt-Grad	4.184(± 0.595)	25.783(± 1.052)	4.183(± 0.592)	25.781(± 1.047)	4.166(± 0.593)	25.781(± 1.050)

	(a) discrepancy (1e-4) and PSNR with $\sigma_b = 10$, $\epsilon = 10$, $\sigma = 12$ on BSD68 (denoising)					
Ord	4.328(± 0.463)	19.334(± 1.017)	4.299(± 0.465)	19.343(± 0.940)	4.289(± 0.461)	19.341(± 0.955)
Jcb	4.039(± 0.440)	19.462(± 1.068)	4.022(± 0.379)	19.472(± 1.029)	4.017(± 0.387)	19.478(± 0.986)
Adv	3.614(± 0.317)	19.561(± 1.102)	3.590(± 0.312)	19.593(± 1.089)	3.589(± 0.320)	19.598(± 1.082)
Smt-Adv	4.030(± 0.560)	20.081(± 1.279)	3.995(± 0.589)	20.163(± 1.233)	3.988(± 0.563)	20.161(± 1.223)
Smt-Grad	2.929(± 0.404)	19.904(± 1.114)	2.917(± 0.407)	19.911(± 1.105)	2.910(± 0.356)	19.943(± 1.110)

	(b) discrepancy (1e-3) and PSNR with $m/n = 0.1$, $\epsilon = 5$, $\sigma = 19$ on Set11 (compressed sensing)					
--	---	--	--	--	--	--

Table 4: Estimated certified discrepancy and PSNR (\pm std.) with different testing samples (S) and the corresponding testing time per data sample (t) in two selected scenarios (a) and (b), where the corresponding best choices of σ (12 and 19) are used.

514 C Details of Algorithms for Randomized Smoothing at Test Time

515 Similar to [4], Algorithm 3 below samples S random realizations of $\delta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ and computes
 516 the corresponding S discrepancies $\ell(f_\theta(\mathbf{y}_{\mathbf{x}+\xi^*} + \delta), f_\theta(\mathbf{y}_{\mathbf{x}+\xi^*}))$. Note that these can be computed
 517 because $\mathbf{y}_{\mathbf{x}+\xi^*}$ is given as an input and $\mathbf{y}_{\mathbf{x}+\xi^*+\delta} = \mathbf{y}_{\mathbf{x}+\xi^*} + \mathbf{A}\delta$ (with known \mathbf{A} and δ). The
 518 sequence of discrepancies is sorted, with the $\lfloor \frac{S}{2} \rfloor$ -th item being used to approximate the median
 519 smoothed discrepancy $D_{0.5}(\mathbf{x}, \xi)$ and the $\lfloor \text{EMPIRICAL}(c, \sigma, S, \epsilon) \rfloor$ -th item to approximate its
 520 certified upper bound $D_{\Phi(\frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0})$.

521 The functionality of EMPIRICAL is described in Algorithm 4. This algorithm is essentially taken
 522 from [7], and despite its rather technical description, is conceptually simple: Given a confidence level
 523 c , perform a binary search to find a value K such that the probability of the event $\text{Binomial}(S, p) \leq K$
 524 is roughly c . Here p is set to $\Phi(\frac{\epsilon}{\sigma})$ in accordance with Corollary 3.2. The idea is then that if we sort
 525 S binomial samples and take the K -th one, we can be confident (up to the confidence level c) that it
 526 is bounded by the actual p -quantile we are interested in. When performing smoothing, we are not
 527 directly using binomial samples, but the indicator variable of whether $\ell(f_\theta(\mathbf{y}_{\mathbf{x}+\xi} + \delta), f_\theta(\mathbf{y}_{\mathbf{x}+\xi}))$ is
 528 below a fixed value (see (6)) indeed follows a binomial distribution. It is worth noting that when
 529 the ratio $\frac{\epsilon}{\sigma}$ is large (e.g., greater than 2), the value of p approaches 1, causing the EMPIRICAL
 530 function to return null. To ensure reliable results, in this paper, we carefully choose ϵ and σ values
 531 that maintain a sufficiently small $\frac{\epsilon}{\sigma}$ ratio to avoid this. Additionally, we adopt the choice of the Monte
 532 Carlo sample size $S = 1000$ and confidence level $c = 0.99$ through out the paper.

Algorithm 3 Find median and certified discrepancy values and corresponding estimates

input Perturbed testing pair $(\mathbf{x}, \mathbf{y}_{\mathbf{x}+\xi^*})$, where ξ^* is the adversarial noise produced by PGD method (but unknown to the algorithm); Estimator f_θ and measurement matrix \mathbf{A} ; Smoothing noise level σ and noise sample count S ; Confidence level c ; Adversarial noise ℓ_2 radius ϵ ; Function $\text{EMPIRICAL}(\cdot)$ from Algorithm 4;
initiate Sorting function $\text{sort}()$; Empty set $\text{discrepancy_set} = \{\}$;

for $1 \leq k \leq S$ **do**
 $\delta_k \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
 $\text{discrepancy_set} \xleftarrow{\text{add}} \ell(f_\theta(\mathbf{y}_{\mathbf{x}+\xi^*+\delta_k}), f_\theta(\mathbf{y}_{\mathbf{x}+\xi^*}))$

end for
 $\text{discrepancy_set_sorted} = \text{sort}(\text{discrepancy_set})$
Set $\hat{D}_{0.5}(\mathbf{x}, \xi) = \text{discrepancy_set_sorted}_{[0.5S]}$
if $\text{EMPIRICAL}(c, \sigma, S, \epsilon)$ returns non-null **then**
 Set $\hat{D}_{\Phi(\frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0}) = \text{discrepancy_set_sorted}_{[\text{EMPIRICAL}(c, \sigma, S, \epsilon)]}$ (Algorithm 4)
 return the estimated median discrepancy $\hat{D}_{0.5}(\mathbf{x}, \xi^*)$, the estimated certificate discrepancy $\hat{D}_{\Phi(\frac{\epsilon}{\sigma})}(\mathbf{x}, \mathbf{0})$, and the two corresponding estimates $\hat{\mathbf{x}} = f_\theta(\mathbf{y}_{\mathbf{x}+\xi^*+\delta_k})$, where $k = [0.5S]$ and $k = [\text{EMPIRICAL}(c, \sigma, S, \epsilon)]$
else
 return median discrepancy and corresponding estimate as above (certified result is unavailable due to overly large ratio $\frac{\epsilon}{\sigma}$)
end if

Algorithm 4 Empirical order statistics generation

input Confidence level c ; Attack radius ϵ ; Smoothing noise level σ ; noise sample count S ;
define Cumulative distribution function of binomial, $F(S, K, p) = \sum_{i=1}^K \binom{S}{i} p^i (1-p)^{S-i}$;

function $\text{EMPIRICAL}(c, \sigma, S, \epsilon)$
 $p = \Phi(\frac{\epsilon}{\sigma})$
 $\underline{K}, \hat{K} = \lceil S \cdot p \rceil, S$
while $\hat{K} - \underline{K} < 1$ **do**
 $\dot{K} = \lfloor \frac{(\hat{K} - \underline{K})}{2} \rfloor$
 if $F(S, \dot{K}, p) > c$ **then**
 $\hat{K} = \dot{K}$
 else
 $\hat{K} = \dot{K}$
 end if
end while
if $\hat{K} < S$ **then**
 return \hat{K}
else
 return null
end if
end function

533 **D Supplementary Experimental Results**

534 Here we provide two additional sets of visual comparisons, Figures 3 and 4, and once again observe
535 that applying randomized smoothing at test time leads to visible visual enhancements in both
536 denoising and compressed sensing tasks across a variety of datasets.

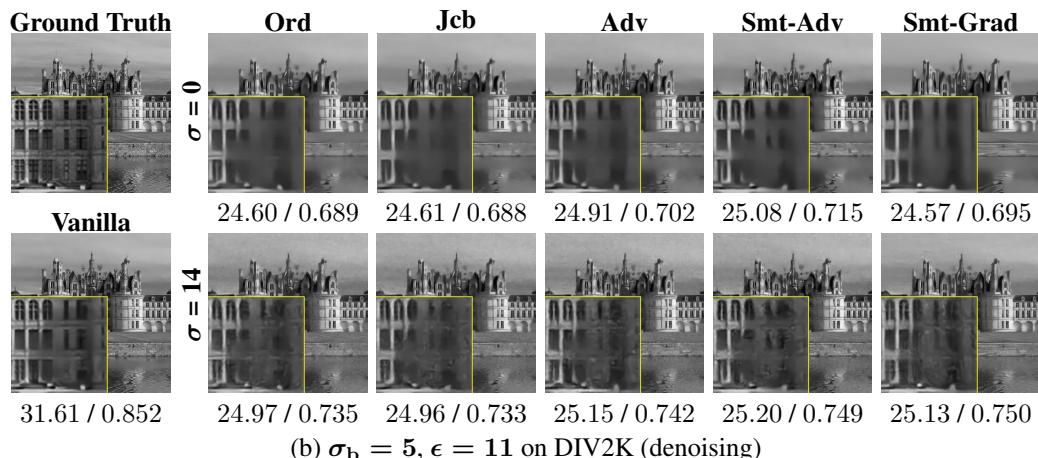
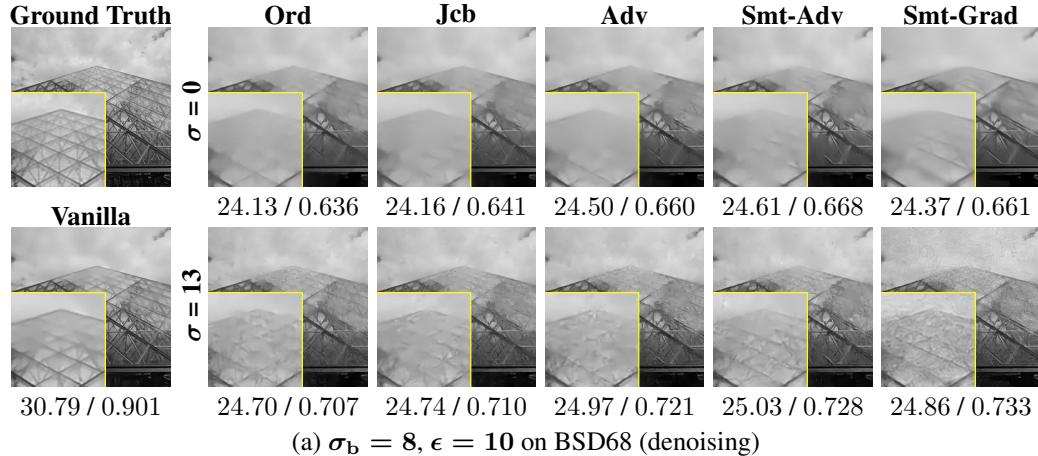
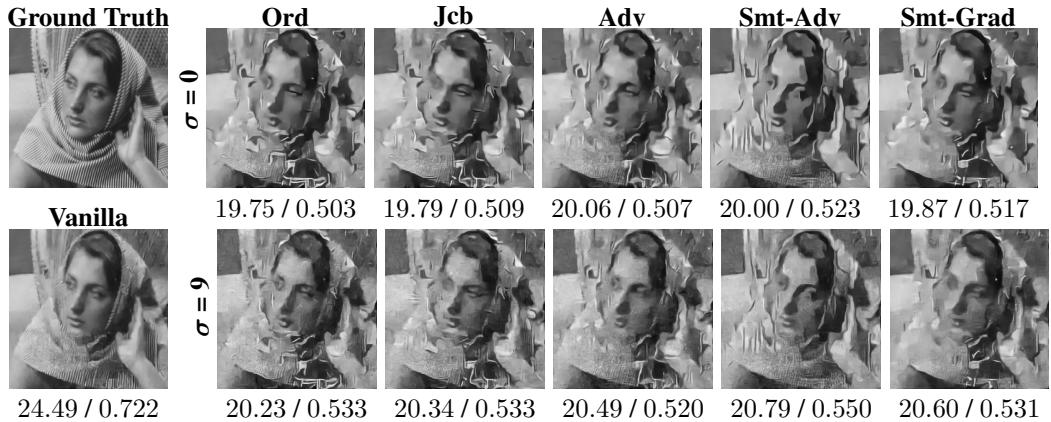
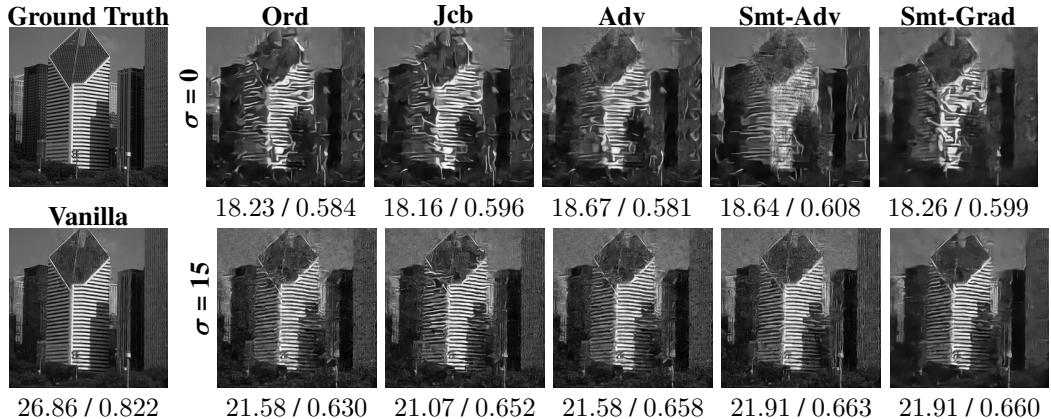


Figure 3: Additional estimations (with the performance shown in the form of PSNR / SSIM) from the perturbed measurements on BSD68 and DIV2K in denoising task.



(a) $m/n = 0.1, \epsilon = 3$ on Set11 (compressed sensing)



(b) $m/n = 0.2, \epsilon = 4$ on ImageNet (compressed sensing)

Figure 4: Additional estimations (with the performance shown in the form of PSNR / SSIM) from the perturbed measurements on Set11 and ImageNet in compressed sensing task.