
DensePure: Understanding Diffusion Models towards Adversarial Robustness

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Diffusion models have been recently employed to improve certified robustness
2 through the process of denoising. However, the theoretical understanding of why
3 diffusion models are able to improve the certified robustness is still lacking, pre-
4 venting from further improvement. In this study, we close this gap by analyzing
5 the fundamental properties of diffusion models and establishing the conditions
6 under which they can enhance certified robustness. This deeper understanding al-
7 lows us to propose a new method DensePure, designed to improve the certified
8 robustness of a pretrained model (i.e. classifier). Given an (adversarial) input,
9 DensePure consists of multiple runs of denoising via the reverse process of the
10 diffusion model (with different random seeds) to get multiple reversed samples,
11 which are then passed through the classifier, followed by majority voting of in-
12 ferred labels to make the final prediction. This design of using multiple runs of
13 denoising is informed by our theoretical analysis of the conditional distribution of
14 the reversed sample. Specifically, when the *data* density of a clean sample is high,
15 its conditional density under the reverse process in a diffusion model is also high;
16 thus sampling from the latter conditional distribution can purify the adversarial
17 example and return the corresponding clean sample with a high probability. By
18 using the highest density point in the conditional distribution as the reversed sam-
19 ple, we identify the robust region of a given instance under the diffusion model’s
20 reverse process. We show that this robust region is a union of multiple convex sets,
21 and is potentially much larger than the robust regions identified in previous works.
22 In practice, DensePure can approximate the label of the high density region in
23 the conditional distribution so that it can enhance certified robustness. We conduct
24 extensive experiments to demonstrate the effectiveness of DensePure by evaluat-
25 ing its certified robustness given a standard model via randomized smoothing. We
26 show that DensePure is consistently better than existing methods on ImageNet,
27 with 7% improvement on average.

28 1 Introduction

29 Diffusion models have been shown to be a powerful image generation tool (Ho et al., 2020; Song
30 et al., 2021b) owing to their iterative diffusion and denoising processes. These models have achieved
31 state-of-the-art performance on sample quality (Dhariwal & Nichol, 2021; Vahdat et al., 2021) as
32 well as effective mode coverage (Song et al., 2021a). A diffusion model usually consists of two
33 processes: (i) a forward diffusion process that converts data to noise by gradually adding noise to
34 the input, and (ii) a reverse generative process that starts from noise and generates data by denoising
35 one step at a time (Song et al., 2021b).

36 Given the natural denoising property of diffusion models, *empirical* studies have leveraged them to
37 perform adversarial purification (Nie et al., 2022; Wu et al., 2022; Carlini et al., 2022). For instance,

38 Nie et al. (2022) introduce a diffusion model based purification model *DiffPure*. They empirically
 39 show that by carefully choosing the amount of Gaussian noises added during the diffusion process,
 40 adversarial perturbations can be removed while preserving the true label semantics. Despite the
 41 significant empirical results, there is no provable guarantee of the achieved robustness. Carlini et al.
 42 (2022) instantiate the randomized smoothing approach with the diffusion model to offer a *provable*
 43 *guarantee* of model robustness against L_2 -norm bounded adversarial example. However, they do
 44 not provide a theoretical understanding of why and how the diffusion models contribute to such
 45 nontrivial certified robustness.

46 **Our Approach.** We theoretically analyze the fundamental properties of diffusion models to under-
 47 stand why and how it enhances certified robustness. This deeper understanding allows us to propose
 48 a new method **DensePure** to improve the certified robustness of any given classifier by more ef-
 49 fectively using the diffusion model. It consists of a pretrained diffusion model and a pretrained
 50 classifier. **DensePure** incorporates two steps: (i) using the reverse process of the diffusion model
 51 to obtain a sample of the posterior data distribution conditioned on the adversarial input; and (ii)
 52 repeating the reverse process multiple times with different random seeds to approximate the label
 53 of high density region in the conditional distribution via a majority vote. In particular, given an ad-
 54 versarial input, we repeatedly feed it into the reverse process of the diffusion model to get multiple
 55 reversed examples and feed them into the classifier to get their labels. We then apply the *majority*
 56 *vote* on the set of labels to get the final predicted label.

57 **DensePure** is inspired by our theoretical analysis, where we show that the diffusion model reverse
 58 process provides a conditional distribution of the reversed sample given an adversarial input, and
 59 sampling from this conditional distribution enhances the certified robustness. Specifically, we prove
 60 that when the data density of clean samples is high, it is a sufficient condition for the conditional
 61 density of the reversed samples to be also high. Therefore, in **DensePure**, samples from the condi-
 62 tional distribution can recover the ground-truth labels with a high probability.

63 For the convenience of understanding and rigorous analysis, we use the highest density point in the
 64 conditional distribution as the deterministic reversed sample for the classifier prediction. We show
 65 that the robust region for a given sample under the diffusion model’s reverse process is the union of
 66 multiple convex sets, each surrounding a region around the ground-truth label. Compared with the
 67 robust region of previous work (Cohen et al., 2019), which only focuses on the neighborhood of *one*
 68 region with the ground-truth label, such union of multiple convex sets has the potential to provide
 69 a much larger robust region. Moreover, the characterization implies that the size of robust regions
 70 is affected by the relative density and the distance between data regions with the ground-truth label
 71 and those with other labels.

72 We conduct extensive experiments on ImageNet and CIFAR-10 datasets under different settings to
 73 evaluate the certifiable robustness of **DensePure**. In particular, we follow the setting from Carlini
 74 et al. (2022) and rely on randomized smoothing to certify robustness to adversarial perturbations
 75 bounded in the \mathcal{L}_2 -norm. We show that **DensePure** achieves the new state-of-the-art *certified*
 76 robustness on the clean model without tuning any model parameters (off-the-shelf). On ImageNet,
 77 it achieves a consistently higher certified accuracy than the existing methods among every σ at every
 78 radius ϵ , 7% improvement on average.

79 2 Preliminaries and Backgrounds

80 **Continuous-Time Diffusion Model.** The diffusion model has two components: the *diffusion pro-*
 81 *cess* followed by the *reverse process*. Given an input random variable $\mathbf{x}_0 \sim p$, the diffusion pro-
 82 cess adds isotropic Gaussian noises to the data so that the diffused random variable at time t is
 83 $\mathbf{x}_t = \sqrt{\alpha_t}(\mathbf{x}_0 + \epsilon_t)$, s.t., $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \sigma_t^2 \mathbf{I})$, and $\sigma_t^2 = (1 - \alpha_t)/\alpha_t$, and we denote $\mathbf{x}_t \sim p_t$. The
 84 forward diffusion process can also be defined by the stochastic differential equation

$$d\mathbf{x} = h(\mathbf{x}, t)dt + g(t)d\mathbf{w}, \quad (\text{SDE})$$

85 where $\mathbf{x}_0 \sim p$, $h : \mathbb{R}^d \times \mathbb{R} \mapsto \mathbb{R}^d$ is the drift coefficient, $g : \mathbb{R} \mapsto \mathbb{R}$ is the diffusion coefficient,
 86 and $\mathbf{w}(t) \in \mathbb{R}^n$ is the standard Wiener process.

87 Under mild conditions C.1, the reverse process exists and removes the added noise by solving the
 88 reverse-time SDE (Anderson, 1982)

$$d\hat{\mathbf{x}} = [h(\hat{\mathbf{x}}, t) - g(t)^2 \nabla_{\hat{\mathbf{x}}} \log p_t(\hat{\mathbf{x}})]dt + g(t)d\bar{\mathbf{w}}, \quad (\text{reverse-SDE})$$

89 where dt is an infinitesimal reverse time step, and $\bar{w}(t)$ is a reverse-time standard Wiener process.
90 In our context, we use the conventions of VP-SDE (Song et al., 2021b) where $h(\mathbf{x}; t) := -\frac{1}{2}\gamma(t)x$
91 and $g(t) := \sqrt{\gamma(t)}$ with $\gamma(t)$ positive and continuous over $[0, 1]$, such that $x(t) = \sqrt{\alpha_t}x(0) +$
92 $\sqrt{1 - \alpha_t}\epsilon$ where $\alpha_t = e^{-\int_0^t \gamma(s)ds}$ and $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. We use $\{\mathbf{x}_t\}_{t \in [0, 1]}$ and $\{\hat{\mathbf{x}}_t\}_{t \in [0, 1]}$ to denote
93 the diffusion process and the reverse process generated by SDE and reverse-SDE respectively, which
94 follow the same distribution.
95 The formulations of Discrete-Time Diffusion Model (or DDPM (Ho et al., 2020)) and Randomized
96 Smoothing are in the appendix.

97 3 Theoretical Analysis

98 In this section, we theoretically analyze why and how the diffusion model can enhance the robustness
99 of a given classifier. We will analyze directly on SDE and reverse-SDE as they generate the same
100 stochastic processes $\{\mathbf{x}_t\}_{t \in [0, T]}$ and the literature works establish an approximation on reverse-
101 SDE (Song et al., 2021b; Ho et al., 2020).

102 We first show that given a diffusion model, solving reverse-SDE will generate a conditional distribu-
103 tion based on the scaled adversarial sample, which will have high density on data region with high
104 *data* density and near to the adversarial sample in Theorem 3.1. See detailed conditions in C.1.

105 **Theorem 3.1.** *Under conditions C.1, solving equation reverse-SDE starting from time t and sample*
106 *$\mathbf{x}_{a,t} = \sqrt{\alpha_t}\mathbf{x}_a$ will generate a reversed random variable $\hat{\mathbf{x}}_0$ with density $\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) \propto$*
107 *$p(\mathbf{x}) \cdot \frac{1}{\sqrt{(2\pi\sigma_t^2)^n}} \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_a\|_2^2}{2\sigma_t^2}\right)$, where p is the data distribution, $\sigma_t^2 = \frac{1 - \alpha_t}{\alpha_t}$ is the variance of*
108 *Gaussian noise added at time t in the diffusion process.*

109 *Proof.* (sketch) Under conditions C.1, we know $\{\mathbf{x}_t\}_{t \in [0, 1]}$ and $\{\hat{\mathbf{x}}_t\}_{t \in [0, 1]}$ follow the same distri-
110 bution, and then the rest proof follows Bayes' Rule. \square

111 Please see the full proofs of this and the following theorems in Appendix C.3.

112 **Remark 1.** *Note that $\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) > 0$ if and only if $p(\mathbf{x}) > 0$, thus the generated reverse*
113 *sample will be on the data region where we train classifiers.*

114 In Theorem 3.1, the conditional density $\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$ is high only if both $p(\mathbf{x})$ and the
115 Gaussian term have high values, i.e., \mathbf{x} has high *data* density and is close to the adversarial sample
116 \mathbf{x}_a . The latter condition is reasonable since adversarial perturbations are typically bounded due to
117 budget constraints. Then, the above argument implies that a reversed sample will have the ground-
118 truth label with a high probability if data region with the ground-truth label has high enough *data*
119 density.

120 For the convenience of theoretical analysis and understanding, we take the point with high-
121 est conditional density $\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$ as the reversed sample, defined as $\mathcal{P}(\mathbf{x}_a; t) :=$
122 $\arg \max_{\mathbf{x}} \mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$. $\mathcal{P}(\mathbf{x}_a; t)$ is a representative of the high density data region in
123 the conditional distribution and $\mathcal{P}(\cdot; t)$ is a deterministic purification model. In the following, we
124 characterize the robust region for data region with ground-truth label under $\mathbb{P}(\cdot; t)$. The robust re-
125 gion and the robust radius for a general deterministic purification model given a classifier are defined
126 below.

127 **Definition 3.2** (Robust Region and Robust Radius). *Given a classifier f and a point \mathbf{x}_0 , let*
128 *$\mathcal{G}(\mathbf{x}_0) := \{\mathbf{x} : f(\mathbf{x}) = f(\mathbf{x}_0)\}$ be the data region where samples have the same label as \mathbf{x}_0 .*
129 *Then given a deterministic purification model $\mathcal{P}(\cdot; \psi)$ with parameter ψ , we define the robust re-*
130 *gion of $\mathcal{G}(\mathbf{x}_0)$ under \mathcal{P} and f as $\mathcal{D}_{\mathcal{P}}^f(\mathcal{G}(\mathbf{x}_0); \psi) := \{\mathbf{x} : f(\mathcal{P}(\mathbf{x}; \psi)) = f(\mathbf{x}_0)\}$, i.e., the set of \mathbf{x}*
131 *such that purified sample $\mathcal{P}(\mathbf{x}; \psi)$ has the same label as \mathbf{x}_0 under f . Further, we define the robust*
132 *radius of \mathbf{x}_0 as $r_{\mathcal{P}}^f(\mathbf{x}_0; \psi) := \max\left\{r : \mathbf{x}_0 + r\mathbf{u} \in \mathcal{D}_{\mathcal{P}}^f(\mathbf{x}_0; \psi), \forall \|\mathbf{u}\|_2 \leq 1\right\}$, i.e., the radius of*
133 *maximum inclined ball of $\mathcal{D}_{\mathcal{P}}^f(\mathbf{x}_0; \psi)$ centered around \mathbf{x}_0 . We will omit \mathcal{P} and f when it is clear*
134 *from the context and write $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); \psi)$ and $r(\mathbf{x}_0; \psi)$ instead.*

135 **Remark 2.** *In Definition 3.2, the robust region (resp. radius) is defined for each class (resp. point).*
136 *When using the point with highest $\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$ as the reversed sample, $\psi := t$.*

137 Now given a sample \mathbf{x}_0 with ground-truth label, we are ready to characterize the robust region
 138 $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); \psi)$ under purification model $\mathcal{P}(\cdot; t)$ and classifier f . Intuitively, if the adversarial sample
 139 \mathbf{x}_a is near to \mathbf{x}_0 (in Euclidean distance), \mathbf{x}_a keeps the same label semantics of \mathbf{x}_0 and so as the
 140 purified sample $\mathcal{P}(\mathbf{x}_a; t)$, which implies that $f(\mathcal{P}(\mathbf{x}_a; \psi)) = f(\mathbf{x}_0)$. However, the condition that
 141 \mathbf{x}_a is near to \mathbf{x}_0 is sufficient but not necessary since we can still achieve $f(\mathcal{P}(\mathbf{x}_a; \psi)) = f(\mathbf{x}_0)$
 142 if \mathbf{x}_a is near to any sample $\tilde{\mathbf{x}}_0$ with $f(\mathcal{P}(\tilde{\mathbf{x}}_0; \psi)) = f(\mathbf{x}_0)$. In the following, we will show that
 143 the robust region $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); \psi)$ is the union of the convex robust sub-regions surrounding every $\tilde{\mathbf{x}}_0$
 144 with the same label as \mathbf{x}_0 . The following theorem characterizes the convex robust sub-region and
 145 robust region respectively.

146 **Theorem 3.3.** *Under conditions C.1 and classifier f , let \mathbf{x}_0 be the sample with ground-truth label
 147 and \mathbf{x}_a be the adversarial sample, then (i) the purified sample $\mathcal{P}(\mathbf{x}_a; t)$ will have the ground-truth
 148 label if \mathbf{x}_a falls into the following convex set,*

$$\mathcal{D}_{\text{sub}}(\mathbf{x}_0; t) := \bigcap_{\{\mathbf{x}'_0: f(\mathbf{x}'_0) \neq f(\mathbf{x}_0)\}} \left\{ \mathbf{x}_a : (\mathbf{x}_a - \mathbf{x}_0)^\top (\mathbf{x}'_0 - \mathbf{x}_0) < \sigma_t^2 \log \left(\frac{p(\mathbf{x}_0)}{p(\mathbf{x}'_0)} \right) + \frac{\|\mathbf{x}'_0 - \mathbf{x}_0\|_2^2}{2} \right\},$$

149 *and further, (ii) the purified sample $\mathcal{P}(\mathbf{x}_a; t)$ will have the ground-truth label if and only if \mathbf{x}_a falls
 150 into the following set, $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t) := \bigcup_{\tilde{\mathbf{x}}_0: f(\tilde{\mathbf{x}}_0) = f(\mathbf{x}_0)} \mathcal{D}_{\text{sub}}(\tilde{\mathbf{x}}_0; t)$. In other words, $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$
 151 is the robust region for data region $\mathcal{G}(\mathbf{x}_0)$ under $\mathcal{P}(\cdot; t)$ and f .*

152 *Proof.* (sketch) (i). Each convex half-space defined by the inequality corresponds to a \mathbf{x}'_0 such that
 153 $f(\mathbf{x}'_0) \neq f(\mathbf{x}_0)$ where \mathbf{x}_a within satisfies $\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x}_0 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) > \mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x}'_0 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$. This
 154 implies that $\mathcal{P}(\mathbf{x}_a; t) \neq \mathbf{x}'_0$ and $f(\mathcal{P}(\mathbf{x}_a; \psi)) = f(\mathbf{x}_0)$. The convexity is due to that the intersection
 155 of convex sets is convex. (ii). The ‘‘if’’ follows directly from (i). The ‘‘only if’’ holds because
 156 if $\mathbf{x}_a \notin \mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$, then exists $\tilde{\mathbf{x}}_1$ such that $f(\tilde{\mathbf{x}}_1) \neq f(\mathbf{x}_0)$ and $\mathbb{P}(\hat{\mathbf{x}}_0 = \tilde{\mathbf{x}}_1 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) >$
 157 $\mathbb{P}(\hat{\mathbf{x}}_0 = \tilde{\mathbf{x}}_0 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}), \forall \tilde{\mathbf{x}}_0$ s.t. $f(\tilde{\mathbf{x}}_0) = f(\mathbf{x}_0)$, and thus $f(\mathcal{P}(\mathbf{x}_a; \psi)) \neq f(\mathbf{x}_0)$. \square

158 **Remark 3.** *Theorem 3.3 implies that when data region $\mathcal{G}(\mathbf{x}_0)$ has higher data density and larger
 159 distances to data regions with other labels, it tends to have larger robust region and points in data
 160 region tends to have larger radius.*

161 In the literature, people focus more on the robust radius (lower bound) $r(\mathcal{G}(\mathbf{x}_0); t)$ (Cohen et al.,
 162 2019; Carlini et al., 2022), which can be obtained by finding the maximum inclined ball inside
 163 $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$ centering \mathbf{x}_0 . Note that although $\mathcal{D}_{\text{sub}}(\mathbf{x}_0; t)$ is convex, $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$ is generally
 164 not. Therefore, finding $r(\mathcal{G}(\mathbf{x}_0); t)$ is a non-convex optimization problem. In particular, it can be
 165 formulated into a disjunctive optimization problem with integer indicator variables, which is typi-
 166 cally NP-hard to solve. One alternative could be finding the maximum inclined ball in $\mathcal{D}_{\text{sub}}(\mathbf{x}_0; t)$,
 167 which can be formulated into a convex optimization problem whose optimal value provides a lower
 168 bound for $r(\mathcal{G}(\mathbf{x}_0); t)$. However, $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$ has the potential to provide much larger robustness
 169 radius because it might connect different convex robust sub-regions into one.

170 In practice, we cannot guarantee to establish an exact reverse process like reverse-SDE but instead
 171 try to establish an approximate reverse process to mimic the exact one. As long as the approximate
 172 reverse process is close enough to the exact reverse process, they will generate close enough con-
 173 ditional distributions based on the adversarial sample. Then the density and locations of the data
 174 regions in two conditional distributions will not differ much and so is the robust region for each
 175 data region. We take the score-based diffusion model in Song et al. (2021b) for an example and
 176 demonstrate Theorem 3.4 to bound the KL-divergence between conditional distributions generated
 177 by reverse-SDE and score-based diffusion model. Ho et al. (2020) showed that using variational
 178 inference to fit DDPM is equivalent to optimizing an objective resembling score-based diffusion
 179 model with a specific weighting scheme, so the results can be extended to DDPM.

180 **Theorem 3.4.** *Under score-based diffusion model Song et al. (2021b) and conditions C.1, we have
 181 $D_{\text{KL}}(\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) \| \mathbb{P}(\mathbf{x}_0^\theta = \mathbf{x} | \mathbf{x}_t^\theta = \mathbf{x}_{a,t})) = \mathcal{I}_{\text{SM}}(\theta, t; \lambda(\cdot))$, where $\{\hat{\mathbf{x}}_\tau\}_{\tau \in [0, t]}$ and
 182 $\{\mathbf{x}_\tau^\theta\}_{\tau \in [0, t]}$ are stochastic processes generated by reverse-SDE and score-based diffusion model
 183 respectively, $\mathcal{I}_{\text{SM}}(\theta, t; \lambda(\cdot)) := \frac{1}{2} \int_0^t \mathbb{E}_{p_\tau(\mathbf{x})} \left[\lambda(\tau) \|\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, \tau)\|_2^2 \right] d\tau$, $\mathbf{s}_\theta(\mathbf{x}, \tau)$ is the
 184 score function to approximate $\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x})$, and $\lambda : \mathbb{R} \rightarrow \mathbb{R}$ is any weighting scheme used in the
 185 training score-based diffusion models.*

Method	Off-the-shelf	Certified Accuracy at ϵ (%)								
		CIFAR-10				ImageNet				
		0.25	0.5	0.75	1.0	0.5	1.0	1.5	2.0	3.0
PixelDP (Lecuyer et al., 2019)	✗	(71.0)22.0	(44.0)2.0	-	-	(33.0)16.0	-	-	-	-
RS (Cohen et al., 2019)	✗	(75.0)61.0	(75.0)43.0	(65.0)32.0	(65.0)23.0	(67.0)49.0	(57.0)37.0	(57.0)29.0	(44.0)19.0	(44.0)12.0
SmoothAdv (Salman et al., 2019a)	✗	(82.0)68.0	(76.0)54.0	(68.0)41.0	(64.0)32.0	(63.0)54.0	(56.0)42.0	(56.0)34.0	(41.0)26.0	(41.0)18.0
Consistency (Jeong & Shim, 2020)	✗	(77.8)68.8	(75.8)58.1	(72.9)48.5	(52.3)37.8	(55.0)50.0	(55.0)44.0	(55.0)34.0	(41.0)24.0	(41.0)17.0
MACER (Zhai et al., 2020)	✗	(81.0)71.0	(81.0)59.0	(66.0)46.0	(66.0)38.0	(68.0)57.0	(64.0)43.0	(64.0)31.0	(48.0)25.0	(48.0)14.0
Boosting (Horváth et al., 2021)	✗	(83.4)70.6	(76.8)60.4	(71.6)52.4	(73.0)38.8	(65.0)57.0	(57.0)44.6	(57.0)38.4	(44.0)28.6	(38.0)21.2
SmoothMix (Jeong et al., 2021)	✓	(77.1)67.9	(77.1)57.9	(74.2)47.7	(61.8)37.2	(55.0)50.0	(55.0)43.0	(55.0)38.0	(40.0)26.0	(40.0)17.0
Denosed (Salman et al., 2020)	✓	(72.0)56.0	(62.0)41.0	(62.0)28.0	(44.0)19.0	(60.0)33.0	(38.0)14.0	(38.0)6.0	-	-
Lee (Lee, 2021)	✓	60.0	42.0	28.0	19.0	41.0	24.0	11.0	-	-
Carlini (Carlini et al., 2022)	✓	(88.0)73.8	(88.0)56.2	(88.0)41.6	(74.2)31.0	(77.0)71.0	(74.0)54.0	(74.0)46.0	(59.0)29.0	(59.0)22.0
Ours	✓	(87.6) 76.6	(87.6) 64.6	(87.6)50.4	(73.6)37.4	(80.0) 76.0	(75.0) 62.0	(75.0) 49.0	(61.0) 37.0	(61.0) 26.0

Table 1: Certified accuracy compared with existing works. The certified accuracy at $\epsilon = 0$ for each model is in the parentheses. The certified accuracy for each cell is from the respective papers except Carlini et al. (2022). Our diffusion model and classifier are the same as Carlini et al. (2022), where the off-the-shelf classifier uses ViT-based architectures trained on a large dataset (ImageNet-22k).

186 *Proof.* (sketch) Let μ_t and ν_t be the path measure for reverse processes $\{\hat{\mathbf{x}}_\tau\}_{\tau \in [0,t]}$ and $\{\mathbf{x}_\tau^\theta\}_{\tau \in [0,t]}$
187 respectively based on the $\mathbf{x}_{a,t}$. Under conditions C.1, μ_t and ν_t are uniquely defined and the KL-
188 divergence can be computed via the Girsanov theorem Oksendal (2013). \square

189 **Remark 4.** *Theorem 3.4 shows that if the training loss is smaller, the conditional distributions gen-*
190 *erated by reverse-SDE and score-based diffusion model are closer, and are the same if the training*
191 *loss is zero.*

192 4 DensePure

193 Inspired by the theoretical analysis, we introduce DensePure and show how to calculate its certified
194 robustness radius via the randomized smoothing algorithm.

195 **Framework.** Our framework, DensePure, consists of two components: (1) an off-the-shelf diffu-
196 sion model with reverse process `rev` and (2) an off-the-shelf base classifier f .

197 Given an input \mathbf{x} , we feed it into the reverse process `rev` of the diffusion model to get
198 the reversed sample `rev(x)` and then repeat the above process K times to get K reversed
199 samples $\{\text{rev}(\mathbf{x})_1, \dots, \text{rev}(\mathbf{x})_K\}$. We feed the above K reversed samples into the clas-
200 sifier to get the corresponding prediction $\{f(\text{rev}(\mathbf{x})_1), \dots, f(\text{rev}(\mathbf{x})_K)\}$ and then apply
201 the *majority vote*, termed **MV**, on these predictions to get the final predicted label $\hat{y} =$
202 $\text{MV}(\{f(\text{rev}(\mathbf{x})_1), \dots, f(\text{rev}(\mathbf{x})_K)\}) = \arg \max_c \sum_{i=1}^K \mathbf{1}\{f(\text{rev}(\mathbf{x})_i) = c\}$.

203 Certified Robustness of DensePure with Randomized Smoothing.

204 We show how DensePure can calculate certified robustness of DensePure via RS, which offers
205 robustness guarantees for a model under a L_2 -norm ball. In particular, we follow the similar setting
206 of Carlini et al. (2022) which uses a DDPM-based diffusion model. The details are in the appendix.

207 5 Experiments

208 In this section, we use DensePure to evaluate certified robustness on two standard datasets, CIFAR-
209 10 (Krizhevsky et al., 2009) and ImageNet (Deng et al., 2009).

210 **Experimental settings** We follow the experimental setting from Carlini et al. (2022). Specifically,
211 for CIFAR-10, we use the 50-M unconditional improved diffusion model from Nichol & Dhariwal
212 (2021) as the diffusion model. We select ViT-B/16 model Dosovitskiy et al. (2020) pretrained on
213 ImageNet-21k and finetuned on CIFAR-10 as the classifier, which could achieve 97.9% accuracy
214 on CIFAR-10. For ImageNet, we use the unconditional 256×256 guided diffusion model from
215 Dhariwal & Nichol (2021) as the diffusion model and pretrained BEiT large model (Bao et al., 2021)
216 trained on ImageNet-21k as the classifier, which could achieve 88.6% top-1 accuracy on validation
217 set of ImageNet-1k. We select three different noise levels $\sigma \in \{0.25, 0.5, 1.0\}$ for certification. For
218 the parameters of DensePure, we set $K = 40$ and $b = 10$ except the results in ablation study. The
219 details about the baselines are in the appendix.

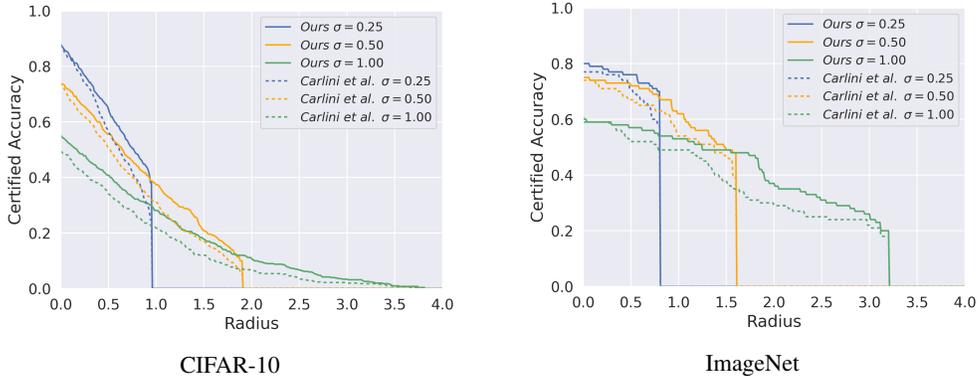


Figure 1: Comparing our method vs Carlini et al. (2022) on CIFAR-10 and ImageNet. The lines represent the certified accuracy with different L_2 perturbation bound with different Gaussian noise $\sigma \in \{0.25, 0.50, 1.00\}$.

220 **Main Results** We compare our results with other baselines. The results are shown in Table 1.

221 For CIFAR-10, comparing with the models which are *carefully* trained with randomized smoothing
 222 techniques in an end-to-end manner (i.e., w/o off-the-shelf classifier), we observe that our method
 223 with the standard off-the-shelf classifier outperforms them at smaller $\epsilon = \{0.25, 0.5\}$ on both
 224 CIFAR-10 and ImageNet datasets while achieves comparable performance at larger $\epsilon = \{0.75, 1.0\}$.
 225 Comparing with the non-diffusion model based methods with off-the-shelf classifier (i.e., De-
 226 noised (Salman et al., 2020) and Lee (Lee, 2021)), both our method and Carlini et al. (2022) are
 227 significantly better than them. These results verify the non-trivial adversarial robustness improve-
 228 ments introduced from the diffusion model. For ImageNet, our method is consistently better than all
 229 priors with a large margin.

230 Since both Carlini et al. (2022) and DensePure use the diffusion model, to better understand the
 231 importance of our design, that approximates the label of the high density region in the conditional
 232 distribution, we compare DensePure with Carlini et al. (2022) in a more fine-grained manner.

233 We show detailed certified robustness of the model among different σ at different radius for CIFAR-
 234 10 in Figure 1-left and for ImageNet in Figure 1-right. We also present our results of certified accu-
 235 racy at different ϵ in Appendix E.3. From these results, we find that our method is still consistently
 236 better at most ϵ (except $\epsilon = 0$) among different σ . The performance margin between ours and Carlini
 237 et al. (2022) will become even larger with a large ϵ . These results further indicate that although the
 238 diffusion model improves model robustness, leveraging the posterior data distribution conditioned
 239 on the input instance (like DensePure) via reverse process instead of using single sample ((Carlini
 240 et al., 2022)) is the key for better robustness. Additionally, we use the off-the-shelf classifiers, which
 241 are the ViT-based architectures trained a larger dataset. In the later ablation study section, we select
 242 the CNN-based architecture wide-ResNet trained on standard dataset from scratch. Our method still
 243 achieves non-trivial robustness.

244 6 Conclusion

245 In this work, we theoretically prove that the diffusion model could purify adversarial examples back
 246 to the corresponding clean sample with high probability, as long as the data density of the cor-
 247 responding clean samples is high enough. Our theoretical analysis characterizes the conditional
 248 distribution of the reversed samples given the adversarial input, generated by the diffusion model
 249 reverse process. Using the highest density point in the conditional distribution as the deterministic
 250 reversed sample, we identify the robust region of a given instance under the diffusion model re-
 251 verse process, which is potentially much larger than previous methods. Our analysis inspires us to
 252 propose an effective pipeline DensePure, for adversarial robustness. We conduct comprehensive
 253 experiments to show the effectiveness of DensePure by evaluating the certified robustness via the
 254 randomized smoothing algorithm. Note that DensePure is an off-the-shelf pipeline that does not
 255 require training a smooth classifier. Our results show that DensePure achieves the new SOTA cer-
 256 tified robustness for perturbation with \mathcal{L}_2 -norm. We hope that our work sheds light on an in-depth
 257 understanding of the diffusion model for adversarial robustness.

258 **References**

- 259 Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Ap-*
260 *plications*, 12(3):313–326, 1982.
- 261 Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint*
262 *arXiv:2106.08254*, 2021.
- 263 Nicholas Carlini, Florian Tramer, J Zico Kolter, et al. (certified!!) adversarial robustness for free!
264 *arXiv preprint arXiv:2206.10550*, 2022.
- 265 Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized
266 smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th*
267 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*
268 *Research*, pp. 1310–1320. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v97/cohen19c.html)
269 [press/v97/cohen19c.html](https://proceedings.mlr.press/v97/cohen19c.html).
- 270 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
271 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
272 pp. 248–255. Ieee, 2009.
- 273 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
274 *in Neural Information Processing Systems*, 34:8780–8794, 2021.
- 275 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
276 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
277 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
278 *arXiv:2010.11929*, 2020.
- 279 Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *Advances*
280 *in Neural Information Processing Systems*, 2019.
- 281 Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi,
282 and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like
283 one. In *International Conference on Learning Representations*, 2020.
- 284 Mitch Hill, Jonathan Craig Mitchell, and Song-Chun Zhu. Stochastic security: Adversarial defense
285 using long-run dynamics of energy-based models. In *International Conference on Learning Rep-*
286 *resentations*, 2021.
- 287 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. URL
288 <https://arxiv.org/abs/2006.11239>.
- 289 Miklós Z Horváth, Mark Niklas Müller, Marc Fischer, and Martin Vechev. Boosting randomized
290 smoothing with variance reduced classifiers. *arXiv preprint arXiv:2106.06946*, 2021.
- 291 Jongheon Jeong and Jinwoo Shin. Consistency regularization for certified robustness of smoothed
292 classifiers. *Advances in Neural Information Processing Systems*, 33:10558–10570, 2020.
- 293 Jongheon Jeong, Sejun Park, Minkyu Kim, Heung-Chang Lee, Do-Guk Kim, and Jinwoo Shin.
294 Smoothmix: Training confidence-calibrated smoothed classifiers for certified robustness. *Ad-*
295 *vances in Neural Information Processing Systems*, 34:30153–30168, 2021.
- 296 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
297 2009.
- 298 Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified
299 robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security*
300 *and Privacy (SP)*, pp. 656–672. IEEE, 2019.
- 301 Kyungmin Lee. Provable defense by denoised smoothing with learned score function. In *ICLR*
302 *Workshop on Security and Safety in Machine Learning Systems*, 2021.

- 303 Matthew Mirman, Timon Gehr, and Martin Vechev. Differentiable abstract interpretation for prov-
304 ably robust neural networks. In *International Conference on Machine Learning*, pp. 3578–3586.
305 PMLR, 2018.
- 306 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
307 In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- 308 Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar.
309 Diffusion models for adversarial purification. In *International Conference on Machine Learning*
310 (*ICML*), 2022.
- 311 Bernt Oksendal. *Stochastic differential equations: an introduction with applications*. Springer
312 Science & Business Media, 2013.
- 313 Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial exam-
314 ples. In *International Conference on Learning Representations*, 2018a.
- 315 Aditi Raghunathan, Jacob Steinhardt, and Percy S Liang. Semidefinite relaxations for certifying
316 robustness to adversarial examples. In *NeurIPS*, 2018b.
- 317 Hadi Salman, Jerry Li, Ilya Razenshteyn, Pengchuan Zhang, Huan Zhang, Sebastien Bubeck, and
318 Greg Yang. Provably robust deep learning via adversarially trained smoothed classifiers. *Ad-
319 vances in Neural Information Processing Systems*, 32, 2019a.
- 320 Hadi Salman, Greg Yang, Huan Zhang, Cho-Jui Hsieh, and Pengchuan Zhang. A convex relax-
321 ation barrier to tight robustness verification of neural networks. *Advances in Neural Information*
322 *Processing Systems*, 32:9835–9846, 2019b.
- 323 Hadi Salman, Mingjie Sun, Greg Yang, Ashish Kapoor, and J Zico Kolter. Denoised smoothing: A
324 provable defense for pretrained classifiers. *Advances in Neural Information Processing Systems*,
325 33:21945–21957, 2020.
- 326 Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-gan: Protecting classifiers against
327 adversarial attacks using generative models. In *International Conference on Learning Representations*,
328 2018.
- 329 Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend:
330 Leveraging generative models to understand and defend against adversarial examples. In *Inter-
331 national Conference on Learning Representations*, 2018.
- 332 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-
333 based diffusion models. *Advances in Neural Information Processing Systems*, 34:1415–1428,
334 2021a.
- 335 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
336 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-
337 ional Conference on Learning Representations*, 2021b.
- 338 Jiachen Sun, Weili Nie, Zhiding Yu, Z Morley Mao, and Chaowei Xiao. Pointdp: Diffusion-
339 driven purification against adversarial attacks on 3d point cloud recognition. *arXiv preprint*
340 *arXiv:2208.09801*, 2022.
- 341 Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space.
342 *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021.
- 343 Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter.
344 Beta-crown: Efficient bound propagation with per-neuron split constraints for neural network
345 robustness verification. *Advances in Neural Information Processing Systems*, 34:29909–29921,
346 2021.
- 347 Quanlin Wu, Hang Ye, and Yuntian Gu. Guided diffusion model for adversarial purification from
348 random noise. *arXiv preprint arXiv:2206.10875*, 2022.
- 349 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint*
350 *arXiv:1605.07146*, 2016.

- 351 Runtian Zhai, Chen Dan, Di He, Huan Zhang, Boqing Gong, Pradeep Ravikumar, Cho-Jui Hsieh,
352 and Liwei Wang. Macer: Attack-free and scalable robust training via maximizing certified radius.
353 *arXiv preprint arXiv:2001.02378*, 2020.
- 354 Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural
355 network robustness certification with general activation functions. In *NeurIPS*, 2018.

356 Appendix

357 Here is the appendix.

358 A Notations

p	data distribution
$\mathbb{P}(A)$	probability of event A
\mathcal{C}^k	set of functions with continuous k -th derivatives
$\mathbf{w}(t)$	standard Wiener Process
$\overline{\mathbf{w}}(t)$	reverse-time standard Wiener Process
$h(\mathbf{x}, t)$	drift coefficient in SDE
$g(t)$	diffusion coefficient in SDE
α_t	scaling coefficient at time t
σ_t^2	variance of added Gaussian noise at time t
$\{\mathbf{x}_t\}_{t \in [0,1]}$	diffusion process generated by SDE
$\{\hat{\mathbf{x}}_t\}_{t \in [0,1]}$	reverse process generated by reverse-SDE
p_t	distribution of \mathbf{x}_t and $\hat{\mathbf{x}}_t$
$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$	diffusion process generated by DDPM
$\{\beta_i\}_{i=1}^N$	pre-defined noise scales in DDPM
ϵ_a	adversarial attack
\mathbf{x}_a	adversarial sample
$\mathbf{x}_{a,t}$	scaled adversarial sample
359 $f(\cdot)$	classifier
$g(\cdot)$	smoothed classifier
$\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$	density of conditional distribution generated by reverse-SDE based on $\mathbf{x}_{a,t}$
$\mathcal{P}(\mathbf{x}_a; t)$	purification model with highest density point
$\mathcal{G}(\mathbf{x}_0)$	data region with the same label as \mathbf{x}_0
$\mathcal{D}_{\mathcal{P}}^f(\mathcal{G}(\mathbf{x}_0); t)$	robust region for $\mathcal{G}(\mathbf{x}_0)$ associated with base classifier f and purification model \mathcal{P}
$r_{\mathcal{P}}^f(\mathbf{x}_0; t)$	robust radius for the point associated with base classifier f and purification model \mathcal{P}
$\mathcal{D}_{sub}(\mathbf{x}_0; t)$	convex robust sub-region
$s_{\theta}(\mathbf{x}, t)$	score function
$\{\mathbf{x}_t^{\theta}\}_{t \in [0,1]}$	reverse process generated by score-based diffusion model
$\mathbb{P}(\mathbf{x}_0^{\theta} = \mathbf{x} \mathbf{x}_t^{\theta} = \mathbf{x}_{a,t})$	density of conditional distribution generated by score-based diffusion model based on $\mathbf{x}_{a,t}$
$\lambda(\tau)$	weighting scheme of training loss for score-based diffusion model
$\mathcal{J}_{SM}(\theta, t; \lambda(\cdot))$	truncated training loss for score-based diffusion model
$\boldsymbol{\mu}_t, \boldsymbol{\nu}_t$	path measure for $\{\hat{\mathbf{x}}_{\tau}\}_{\tau \in [0,t]}$ and $\{\mathbf{x}_{\tau}^{\theta}\}_{\tau \in [0,t]}$ respectively

360 B Related Work

361 Using an off-the-shelf generative model to purify adversarial perturbations has become an important
 362 direction in adversarial defense. Previous works have developed various purification methods based
 363 on different generative models, such as GANs (Samangouei et al., 2018), autoregressive generative
 364 models (Song et al., 2018), and energy-based models (Du & Mordatch, 2019; Grathwohl et al.,
 365 2020; Hill et al., 2021). More recently, as diffusion models (or score-based models) achieve better
 366 generation quality than other generative models (Ho et al., 2020; Dhariwal & Nichol, 2021), many
 367 works consider using diffusion models for adversarial purification (Nie et al., 2022; Wu et al., 2022;
 368 Sun et al., 2022) Although they have found good empirical results in defending against existing
 369 adversarial attacks (Nie et al., 2022), there is no provable guarantee about the robustness about such
 370 methods. On the other hand, certified defenses provide guarantees of robustness (Mirman et al.,
 371 2018; Cohen et al., 2019; Lecuyer et al., 2019; Salman et al., 2020; Horváth et al., 2021; Zhang et al.,
 372 2018; Raghunathan et al., 2018a,b; Salman et al., 2019b; Wang et al., 2021). They provide a lower
 373 bounder of model accuracy under constrained perturbations. Among them, approaches Lecuyer et al.
 374 (2019); Cohen et al. (2019); Salman et al. (2019a); Jeong & Shin (2020); Zhai et al. (2020); Horváth
 375 et al. (2021); Jeong et al. (2021); Salman et al. (2020); Lee (2021); Carlini et al. (2022) based
 376 on randomized smoothing (Cohen et al., 2019) show the great scalability and achieve promising
 377 performance on large network and dataset. The most similar work to us is Carlini et al. (2022), which
 378 uses diffusion models combined with standard classifiers for certified defense. They view diffusion
 379 model as blackbox without having a theoretical understanding of why and how the diffusion models
 380 contribute to such nontrivial certified robustness.

381 C More details about Theoretical analysis

382 C.1 Assumptions

- 383 (i) The data distribution $p \in \mathcal{C}^2$ and $\mathbb{E}_{\mathbf{x} \sim p} [\|\mathbf{x}\|_2^2] < \infty$.
- 384 (ii) $\forall t \in [0, T] : h(\cdot, t) \in \mathcal{C}^1, \exists C > 0, \forall \mathbf{x} \in \mathbb{R}^n, t \in [0, T] : \|h(\mathbf{x}, t)\|_2 \leq C(1 + \|\mathbf{x}\|_2)$.
- 385 (iii) $\exists C > 0, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n : \|h(\mathbf{x}, t) - h(\mathbf{y}, t)\|_2 \leq C\|\mathbf{x} - \mathbf{y}\|_2$.
- 386 (iv) $g \in \mathcal{C}$ and $\forall t \in [0, T], |g(t)| > 0$.
- 387 (v) $\forall t \in [0, T] : s_\theta(\cdot, t) \in \mathcal{C}^1, \exists C > 0, \forall \mathbf{x} \in \mathbb{R}^n, t \in [0, T] : \|s_\theta(\mathbf{x}, t)\|_2 \leq C(1 + \|\mathbf{x}\|_2)$.
- 388 (vi) $\exists C > 0, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n : \|s_\theta(\mathbf{x}, t) - s_\theta(\mathbf{y}, t)\|_2 \leq C\|\mathbf{x} - \mathbf{y}\|_2$.

389 C.2 Background

390 **Discrete-Time Diffusion Model (or DDPM (Ho et al., 2020)).** DDPM constructs a discrete
 391 Markov chain $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N\}$ as the forward process for the training data $\mathbf{x}_0 \sim p$, such
 392 that $\mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1}) = \mathcal{N}(\mathbf{x}_i; \sqrt{1 - \beta_i} \mathbf{x}_{i-1}, \beta_i \mathbf{I})$, where $0 < \beta_1 < \beta_2 < \dots < \beta_N < 1$ are predefined
 393 noise scales such that \mathbf{x}_N approximates the Gaussian white noise. Denote $\bar{\alpha}_i = \prod_{j=1}^i (1 - \beta_j)$, we
 394 have $\mathbb{P}(\mathbf{x}_i | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_i; \sqrt{\bar{\alpha}_i} \mathbf{x}_0, (1 - \bar{\alpha}_i) \mathbf{I})$, i.e., $\mathbf{x}_t(\mathbf{x}_0, \epsilon) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + (1 - \bar{\alpha}_t) \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

395 The reverse process of DDPM learns a reverse direction variational Markov chain $p_\theta(\mathbf{x}_{i-1} | \mathbf{x}_i) =$
 396 $\mathcal{N}(\mathbf{x}_{i-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_i, i), \Sigma_\theta(\mathbf{x}_i, i))$. Ho et al. (2020) defines ϵ_θ as a function approximator to predict
 397 ϵ from \mathbf{x}_i such that $\boldsymbol{\mu}_\theta(\mathbf{x}_i, i) = \frac{1}{\sqrt{1 - \beta_i}} \left(\mathbf{x}_i - \frac{\beta_i}{\sqrt{1 - \bar{\alpha}_i}} \epsilon_\theta(\mathbf{x}_i, i) \right)$. Then the reverse time samples
 398 are generated by $\hat{\mathbf{x}}_{i-1} = \frac{1}{\sqrt{1 - \beta_i}} \left(\hat{\mathbf{x}}_i - \frac{\beta_i}{\sqrt{1 - \bar{\alpha}_i}} \epsilon_{\theta^*}(\hat{\mathbf{x}}_i, i) \right) + \sqrt{\beta_i} \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and the optimal
 399 parameters θ^* are obtained by solving $\theta^* := \arg \min_\theta \mathbb{E}_{\mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + (1 - \bar{\alpha}_i) \epsilon)\|_2^2]$.

400 **Randomized Smoothing.** Randomized smoothing is used to certify the robustness of a given
 401 classifier against L_2 -norm based perturbation. It transfers the classifier f to a smooth version
 402 $g(\mathbf{x}) = \arg \max_c \mathbb{P}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} (f(\mathbf{x} + \epsilon) = c)$, where g is the smooth classifier and σ is a hyper-
 403 parameter of the smooth classifier g , which controls the trade-off between robustness and accuracy.
 404 Cohen et al. (2019) shows that $g(x)$ induces the certifiable robustness for \mathbf{x} under the L_2 -norm with
 405 radius R , where $R = \frac{\sigma}{2} (\Phi^{-1}(p_A) - \Phi^{-1}(p_B))$; p_A and p_B are probability of the most probable
 406 class and ‘‘runner-up’’ class respectively; Φ is the inverse of the standard Gaussian CDF. The p_A and
 407 p_B can be estimated with arbitrarily high confidence via Monte Carlo method (Cohen et al., 2019).

408 **C.3 Theorems and Proofs**

409 **Theorem 3.1.** Under conditions C.1, solving equation reverse-SDE starting from time t and point
 410 $\mathbf{x}_{a,t} = \sqrt{\alpha_t} \mathbf{x}_a$ will generate a reversed random variable $\hat{\mathbf{x}}_0$ with conditional distribution

$$\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) \propto p(\mathbf{x}) \cdot \frac{1}{\sqrt{(2\pi\sigma_t^2)^n}} e^{-\frac{\|\mathbf{x} - \mathbf{x}_a\|_2^2}{2\sigma_t^2}}$$

411 where $\sigma_t^2 = \frac{1-\alpha_t}{\alpha_t}$ is the variance of the Gaussian noise added at timestamp t in the diffusion
 412 process SDE.

413 *Proof.* Under the assumption, we know $\{\mathbf{x}_t\}_{t \in [0,1]}$ and $\{\hat{\mathbf{x}}_t\}_{t \in [0,1]}$ follow the same distribution,
 414 which means

$$\begin{aligned} \mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) &= \frac{\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x}, \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})}{\mathbb{P}(\hat{\mathbf{x}}_t = \mathbf{x}_{a,t})} \\ &= \frac{\mathbb{P}(\mathbf{x}_0 = \mathbf{x}, \mathbf{x}_t = \mathbf{x}_{a,t})}{\mathbb{P}(\mathbf{x}_t = \mathbf{x}_{a,t})} \\ &= \mathbb{P}(\mathbf{x}_0 = \mathbf{x}) \frac{\mathbb{P}(\mathbf{x}_t = \mathbf{x}_{a,t} | \mathbf{x}_0 = \mathbf{x})}{\mathbb{P}(\mathbf{x}_t = \mathbf{x}_{a,t})} \\ &\propto \mathbb{P}(\mathbf{x}_0 = \mathbf{x}) \frac{1}{\sqrt{(2\pi\sigma_t^2)^n}} e^{-\frac{\|\mathbf{x} - \mathbf{x}_a\|_2^2}{2\sigma_t^2}} \\ &= p(\mathbf{x}) \cdot \frac{1}{\sqrt{(2\pi\sigma_t^2)^n}} e^{-\frac{\|\mathbf{x} - \mathbf{x}_a\|_2^2}{2\sigma_t^2}} \end{aligned}$$

415 where the third equation is due to the chain rule of probability and the last equation is a result of the
 416 diffusion process. \square

417 **Theorem 3.3.** Under conditions C.1 and classifier f , let \mathbf{x}_0 be the sample with ground-truth label
 418 and \mathbf{x}_a be the adversarial sample, then (i) the purified sample $\mathcal{P}(\mathbf{x}_a; t)$ will have the ground-truth
 419 label if \mathbf{x}_a falls into the following convex set,

$$\mathcal{D}_{sub}(\mathbf{x}_0; t) := \bigcap_{\{\mathbf{x}'_0: f(\mathbf{x}'_0) \neq f(\mathbf{x}_0)\}} \left\{ \mathbf{x}_a : (\mathbf{x}_a - \mathbf{x}_0)^\top (\mathbf{x}'_0 - \mathbf{x}_0) < \sigma_t^2 \log \left(\frac{p(\mathbf{x}_0)}{p(\mathbf{x}'_0)} \right) + \frac{\|\mathbf{x}'_0 - \mathbf{x}_0\|_2^2}{2} \right\},$$

420 and further, (ii) the purified sample $\mathcal{P}(\mathbf{x}_a; t)$ will have the ground-truth label if and only if \mathbf{x}_a falls
 421 into the following set, $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t) := \bigcup_{\tilde{\mathbf{x}}_0: f(\tilde{\mathbf{x}}_0) = f(\mathbf{x}_0)} \mathcal{D}_{sub}(\tilde{\mathbf{x}}_0; t)$. In other words, $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$
 422 is the robust region for data region $\mathcal{G}(\mathbf{x}_0)$ under $\mathcal{P}(\cdot; t)$ and f .

423 *Proof.* We start with part (i).

424 The main idea is to prove that a point \mathbf{x}'_0 such that $f(\mathbf{x}'_0) \neq f(\mathbf{x}_0)$ should have lower density than
 425 \mathbf{x}_0 in the conditional distribution in Theorem 3.1 so that $\mathcal{P}(\mathbf{x}_a; t)$ cannot be \mathbf{x}'_0 . In other words, we
 426 should have

$$\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x}_0 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) > \mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x}'_0 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}).$$

427 By Theorem 3.1, this is equivalent to

$$\begin{aligned} p(\mathbf{x}_0) \cdot \frac{1}{\sqrt{(2\pi\sigma_t^2)^n}} e^{-\frac{\|\mathbf{x}_0 - \mathbf{x}_a\|_2^2}{2\sigma_t^2}} &> p(\mathbf{x}'_0) \cdot \frac{1}{\sqrt{(2\pi\sigma_t^2)^n}} e^{-\frac{\|\mathbf{x}'_0 - \mathbf{x}_a\|_2^2}{2\sigma_t^2}} \\ \Leftrightarrow \log \left(\frac{p(\mathbf{x}_0)}{p(\mathbf{x}'_0)} \right) &> \frac{1}{2\sigma_t^2} (\|\mathbf{x}_0 - \mathbf{x}_a\|_2^2 - \|\mathbf{x}'_0 - \mathbf{x}_a\|_2^2) \\ \Leftrightarrow \log \left(\frac{p(\mathbf{x}_0)}{p(\mathbf{x}'_0)} \right) &> \frac{1}{2\sigma_t^2} (\|\mathbf{x}_0 - \mathbf{x}_a\|_2^2 - \|\mathbf{x}'_0 - \mathbf{x}_0 + \mathbf{x}_0 - \mathbf{x}_a\|_2^2) \\ \Leftrightarrow \log \left(\frac{p(\mathbf{x}_0)}{p(\mathbf{x}'_0)} \right) &> \frac{1}{2\sigma_t^2} (2(\mathbf{x}_a - \mathbf{x}_0)^\top (\mathbf{x}'_0 - \mathbf{x}_0) - \|\mathbf{x}'_0 - \mathbf{x}_0\|_2^2). \end{aligned}$$

428 Re-organizing the above inequality, we obtain

$$(\mathbf{x}_a - \mathbf{x}_0)^\top (\mathbf{x}'_0 - \mathbf{x}_0) < \sigma_t^2 \log \left(\frac{p(\mathbf{x}_0)}{p(\mathbf{x}'_0)} \right) + \frac{1}{2} \|\mathbf{x}'_0 - \mathbf{x}_0\|_2^2.$$

429 Note that the order of \mathbf{x}_a is at most one in every term of the above inequality, so the inequality
 430 actually defines a half-space in \mathbb{R}^n for every $(\mathbf{x}_0, \mathbf{x}'_0)$ pair. Further, we have to satisfy the inequality
 431 for every \mathbf{x}'_0 such that $f(\mathbf{x}'_0) \neq f(\mathbf{x}_0)$, therefore, by intersecting over all such half-spaces, we
 432 obtain a convex $\mathcal{D}_{\text{sub}}(\mathbf{x}_0; t)$.

433 Then we prove part (ii).

434 On the one hand, if $\mathbf{x}_a \in \mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$, then there exists one $\tilde{\mathbf{x}}_0$ such that $f(\tilde{\mathbf{x}}_0) = f(\mathbf{x}_0)$ and
 435 $\mathbf{x}_a \in \mathcal{D}_{\text{sub}}(\tilde{\mathbf{x}}_0; t)$. By part (i), $\tilde{\mathbf{x}}_0$ has higher probability than all other points with different labels
 436 from \mathbf{x}_0 in the conditional distribution $\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$ characterized by Theorem 3.1.
 437 Therefore, $\mathcal{P}(\mathbf{x}_a; t)$ should have the same label as \mathbf{x}_0 . On the other hand, if $\mathbf{x}_a \notin \mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$,
 438 then there is a point $\tilde{\mathbf{x}}_1$ with different label from \mathbf{x}_0 such that for any $\tilde{\mathbf{x}}_0$ with the same label as \mathbf{x}_0 ,
 439 $\mathbb{P}(\hat{\mathbf{x}}_0 = \tilde{\mathbf{x}}_1 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) > \mathbb{P}(\hat{\mathbf{x}}_0 = \tilde{\mathbf{x}}_0 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$. In other words, $\mathcal{P}(\mathbf{x}_a; t)$ would have different
 440 label from \mathbf{x}_0 . \square

441 **Theorem 3.4.** *Under score-based diffusion model Song et al. (2021b) and conditions C.1, we can*
 442 *bound*

$$D_{\text{KL}}(\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) \| \mathbb{P}(\mathbf{x}_0^\theta = \mathbf{x} | \mathbf{x}_t^\theta = \mathbf{x}_{a,t})) = \mathcal{J}_{\text{SM}}(\theta, t; \lambda(\cdot))$$

where $\{\hat{\mathbf{x}}_\tau\}_{\tau \in [0,t]}$ and $\{\mathbf{x}_\tau^\theta\}_{\tau \in [0,t]}$ are stochastic processes generated by reverse-SDE and score-
 based diffusion model respectively,

$$\mathcal{J}_{\text{SM}}(\theta, t; \lambda(\cdot)) := \frac{1}{2} \int_0^t \mathbb{E}_{p_\tau(\mathbf{x})} \left[\lambda(\tau) \|\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, \tau)\|_2^2 \right] d\tau,$$

443 $\mathbf{s}_\theta(\mathbf{x}, \tau)$ is the score function to approximate $\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x})$, and $\lambda : \mathbb{R} \rightarrow \mathbb{R}$ is any weighting scheme
 444 used in the training score-based diffusion models.

445 *Proof.* Similar to proof of (Song et al., 2021a, Theorem 1), let μ_t and ν_t be the path measure for
 446 reverse processes $\{\hat{\mathbf{x}}_\tau\}_{\tau \in [0,t]}$ and $\{\mathbf{x}_\tau^\theta\}_{\tau \in [0,t]}$ respectively based on the scaled adversarial sample
 447 $\mathbf{x}_{a,t}$. Under conditions C.1, the KL-divergence can be computed via the Girsanov theorem Oksendal
 448 (2013):

$$\begin{aligned} & D_{\text{KL}}(\mathbb{P}(\hat{\mathbf{x}}_0 = \mathbf{x} | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) \| \mathbb{P}(\mathbf{x}_0^\theta = \mathbf{x} | \mathbf{x}_t^\theta = \mathbf{x}_{a,t})) \\ &= -\mathbb{E}_{\mu_t} \left[\log \frac{d\nu_t}{d\mu_t} \right] \\ &\stackrel{(i)}{=} \mathbb{E}_{\mu_t} \left[\int_0^t g(\tau) (\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, \tau)) d\bar{\mathbf{w}}_\tau + \frac{1}{2} \int_0^t g(\tau)^2 \|\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, \tau)\|_2^2 d\tau \right] \\ &= \mathbb{E}_{\mu_t} \left[\frac{1}{2} \int_0^t g(\tau)^2 \|\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, \tau)\|_2^2 d\tau \right] \\ &= \frac{1}{2} \int_0^t \mathbb{E}_{p_\tau(\mathbf{x})} \left[g(\tau)^2 \|\nabla_{\mathbf{x}} \log p_\tau(\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}, \tau)\|_2^2 \right] d\tau \\ &= \mathcal{J}_{\text{SM}}(\theta, t; g(\cdot)^2) \end{aligned}$$

449 where (i) is due to Girsanov Theorem and (ii) is due to the martingale property of Itô integrals. \square

450 D More details about DensePure

451 D.1 Pseudo-Code

452 We provide the pseudo code of DensePure in Algo. 1 and Alg. 2

Algorithm 1 DensePure pseudo-code with the highest density point

- 1: Initialization: choose off-the-shelf diffusion model and classifier f , choose $\psi = t$,
 - 2: Input sample $\mathbf{x}_a = \mathbf{x}_0 + \epsilon_a$
 - 3: Compute $\hat{\mathbf{x}}_0 = \mathcal{P}(\mathbf{x}_a; \psi)$
 - 4: $\hat{y} = f(\hat{\mathbf{x}}_0)$
-

Algorithm 2 DensePure pseudo-code with majority vote

- 1: Initialization: choose off-the-shelf diffusion model and classifier f , choose σ
 - 2: Compute $\bar{\alpha}_n = \frac{1}{1+\sigma^2}$, $n = \arg \min_s \left\{ \left| \bar{\alpha}_s - \frac{1}{1+\sigma^2} \right| \mid s \in \{1, 2, \dots, N\} \right\}$
 - 3: Generate input sample $\mathbf{x}_{rs} = \mathbf{x}_0 + \epsilon$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
 - 4: Choose schedule S^b , get $\hat{\mathbf{x}}_0^i \leftarrow \text{rev}(\sqrt{\bar{\alpha}_n} \mathbf{x}_{rs})_i$, $i = 1, 2, \dots, K$ with Fast Sampling
 - 5: $\hat{y} = \mathbf{MV}(\{f(\hat{\mathbf{x}}_0^1), \dots, f(\hat{\mathbf{x}}_0^K)\}) = \arg \max_c \sum_{i=1}^K \mathbf{1}\{f(\hat{\mathbf{x}}_0^i) = c\}$
-

453 D.2 Certified Robustness of DensePure with Randomized Smoothing.

454 We show how DensePure can calculate certified robustness of DensePure via RS, which offers
455 robustness guarantees for a model under a L_2 -norm ball.

456 In particular, we follow the similar setting of Carlini et al. (2022) which uses a DDPM-based diffu-
457 sion model. The details are in the appendix. The overall algorithm contains three steps:

458 (1) Our framework estimates n , the number of steps used for the reverse process of DDPM-based
459 diffusion model. Since Randomized Smoothing (Cohen et al., 2019) adds Gaussian noise ϵ , where
460 $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, to data input \mathbf{x} to get the randomized data input, $\mathbf{x}_{rs} = \mathbf{x} + \epsilon$, we map between
461 the noise required by the randomized example \mathbf{x}_{rs} and the noise required by the diffused data \mathbf{x}_n
462 (i.e., $\mathbf{x}_n \sim \mathcal{N}(\mathbf{x}_n; \sqrt{\bar{\alpha}_n} \mathbf{x}_0, (1 - \bar{\alpha}_n) \mathbf{I})$) with n step diffusion processing so that $\bar{\alpha}_n = \frac{1}{1+\sigma^2}$. In
463 this way, we can compute the corresponding timestep n , where $n = \arg \min_s \{|\bar{\alpha}_s - \frac{1}{1+\sigma^2}| \mid s \in$
464 $\{1, 2, \dots, N\}\}$.

465 (2). Given the above calculated timestep n , we scale \mathbf{x}_{rs} with $\sqrt{\bar{\alpha}_n}$ to obtain the scaled randomized
466 smoothing sample $\sqrt{\bar{\alpha}_n} \mathbf{x}_{rs}$. Then we feed $\sqrt{\bar{\alpha}_n} \mathbf{x}_{rs}$ into the reverse process of the diffusion model
467 by K -times to get the reversed sample set $\{\hat{\mathbf{x}}_0^1, \hat{\mathbf{x}}_0^2, \dots, \hat{\mathbf{x}}_0^i, \dots, \hat{\mathbf{x}}_0^K\}$.

468 (3). We feed the obtained reversed sample set into a standard *off-the-shelf* classifier f to get the
469 corresponding predicted labels $\{f(\hat{\mathbf{x}}_0^1), f(\hat{\mathbf{x}}_0^2), \dots, f(\hat{\mathbf{x}}_0^i), \dots, f(\hat{\mathbf{x}}_0^K)\}$, and apply *majority vote*,
470 denoted $\mathbf{MV}(\dots)$, on these predicted labels to get the final label for \mathbf{x}_{rs} .

471 To calculate the reversed sample, the standard reverse process of DDPM-based models re-
472 quire repeatedly applying a “single-step” operation n times to get the reversed sample $\hat{\mathbf{x}}_0$
473 (i.e., $\hat{\mathbf{x}}_0 = \underbrace{\text{Reverse}(\dots \text{Reverse}(\dots \text{Reverse}(\text{Reverse}(\sqrt{\bar{\alpha}_n} \mathbf{x}_{rs}; n); n-1); \dots); i)}_{n \text{ steps}}(\dots 1)$). Here

474 $\hat{\mathbf{x}}_{i-1} = \text{Reverse}(\hat{\mathbf{x}}_i; i)$ is equivalent to sample $\hat{\mathbf{x}}_{i-1}$ from $\mathcal{N}(\hat{\mathbf{x}}_{i-1}; \boldsymbol{\mu}_\theta(\hat{\mathbf{x}}_i, i), \boldsymbol{\Sigma}_\theta(\hat{\mathbf{x}}_i, i))$, where
475 $\boldsymbol{\mu}_\theta(\hat{\mathbf{x}}_i, i) = \frac{1}{\sqrt{1-\beta_i}} \left(\hat{\mathbf{x}}_i - \frac{\beta_i}{\sqrt{1-\bar{\alpha}_i}} \boldsymbol{\epsilon}_\theta(\hat{\mathbf{x}}_i, i) \right)$ and $\boldsymbol{\Sigma}_\theta := \exp(v \log \beta_i + (1-v) \log \tilde{\beta}_i)$. Here v is a
476 parameter learned by DDPM and $\tilde{\beta}_i = \frac{1-\bar{\alpha}_{i-1}}{1-\bar{\alpha}_i}$.

477 To reduce the time complexity, we use the uniform sub-sampling strategy from Nichol & Dhariwal
478 (2021). We uniformly sample a subsequence with size b from the original N -step the reverse process.

479 In details, we follow the method used in (Nichol & Dhariwal, 2021) and sample a subsequence
480 S^b with b values (i.e., $S^b = \underbrace{\{n, \lfloor n - \frac{n}{b} \rfloor, \dots, 1\}}_b$, where S_i^b is the i -th element in S^b and $S_i^b =$

481 $\lfloor n - \frac{in}{b} \rfloor, \forall i < b$ and $S_b^b = 1$) from the original schedule S (i.e., $S = \underbrace{\{n, n-1, \dots, 1\}}_n$, where

482 $S_i = i$ is the i -th element in S).

Methods	Noise	Certified Accuracy at ϵ (%)				
		0.0	0.25	0.5	0.75	1.0
Carlini (Carlini et al., 2022)	$\sigma = 0.25$	88.0	73.8	56.2	41.6	0.0
	$\sigma = 0.5$	74.2	62.0	50.4	40.2	31.0
	$\sigma = 1.0$	49.4	41.4	34.2	27.8	21.8
Ours	$\sigma = 0.25$	87.6(-0.4)	76.6(+2.8)	64.6(+8.4)	50.4(+8.8)	0.0(+0.0)
	$\sigma = 0.5$	73.6(-0.6)	65.4(+3.4)	55.6(+5.2)	46.0(+5.8)	37.4(+6.4)
	$\sigma = 1.0$	55.0(+5.6)	47.8(+6.4)	40.8(+6.6)	33.0(+5.2)	28.2(+6.4)

Table A: Certified accuracy compared with Carlini et al. (2022) for CIFAR-10 at all σ . The numbers in the bracket are the difference of certified accuracy between two methods. Our diffusion model and classifier are the same as Carlini et al. (2022).

483 Within this context, we adapt the original $\bar{\alpha}$ schedule $\bar{\alpha}^S = \{\bar{\alpha}_1, \dots, \bar{\alpha}_i, \dots, \bar{\alpha}_n\}$ used for single-
484 step to the new schedule $\bar{\alpha}^{S^b} = \{\bar{\alpha}_{S_1^b}, \dots, \bar{\alpha}_{S_j^b}, \dots, \bar{\alpha}_{S_b^b}\}$ (i.e., $\bar{\alpha}_i^{S^b} = \bar{\alpha}_{S_i^b} = \bar{\alpha}_{S_{\lfloor n - \frac{i}{b} \rfloor}}$ is the
485 i -th element in $\bar{\alpha}^{S^b}$). We calculate the corresponding $\beta^{S^b} = \{\beta_1^{S^b}, \beta_2^{S^b}, \dots, \beta_i^{S^b}, \dots, \beta_b^{S^b}\}$ and
486 $\tilde{\beta}^{S^b} = \{\tilde{\beta}_1^{S^b}, \tilde{\beta}_2^{S^b}, \dots, \tilde{\beta}_i^{S^b}, \dots, \tilde{\beta}_b^{S^b}\}$ schedules, where $\beta_{S_i^b} = \beta_i^{S^b} = 1 - \frac{\bar{\alpha}_i^{S^b}}{\bar{\alpha}_{i-1}^{S^b}}$, $\tilde{\beta}_{S_i^b} =$
487 $\tilde{\beta}_i^{S^b} = \frac{1 - \bar{\alpha}_{i-1}^{S^b}}{1 - \bar{\alpha}_i^{S^b}} \beta_{S_i^b}$. With these new schedules, we can use b times reverse steps to calculate
488 $\hat{x}_0 = \underbrace{\text{Reverse}(\dots \text{Reverse}(\text{Reverse}(x_n; S_b^b); S_{b-1}^b); \dots; 1)}_b$. Since $\Sigma_{\theta}(x_{S_i^b}, S_i^b)$ is parameterized
489 as a range between β^{S^b} and $\tilde{\beta}^{S^b}$, it will automatically be rescaled. Thus, $\hat{x}_{S_{i-1}^b} = \text{Reverse}(\hat{x}_{S_i^b}; S_i^b)$
490 is equivalent to sample $x_{S_{i-1}^b}$ from $\mathcal{N}(x_{S_i^b}; \mu_{\theta}(x_{S_i^b}, S_i^b), \Sigma_{\theta}(x_{S_i^b}, S_i^b))$.

491 E More Experimental details and Results

492 E.1 Implementation details

493 We select three different noise levels $\sigma \in \{0.25, 0.5, 1.0\}$ for certification. For the parameters
494 of DensePure, The sampling numbers when computing the certified radius are $n = 100000$ for
495 CIFAR-10 and $n = 10000$ for ImageNet. We evaluate the certified robustness on 500 samples subset
496 of CIFAR-10 testset and 100 samples subset of ImageNet validation set. we set $K = 40$ and $b = 10$
497 except the results in ablation study. The details about the baselines are in the appendix.

498 E.2 Baselines.

499 We select randomized smoothing based methods including PixelDP (Lecuyer et al., 2019), RS (Co-
500 hen et al., 2019), SmoothAdv (Salman et al., 2019a), Consistency (Jeong & Shin, 2020), MACER
501 (Zhai et al., 2020), Boosting (Horváth et al., 2021), SmoothMix (Jeong et al., 2021), Denoised
502 (Salman et al., 2020), Lee (Lee, 2021), Carlini (Carlini et al., 2022) as our baselines. Among them,
503 PixelDP, RS, SmoothAdv, Consistency, MACER, and SmoothMix require training a smooth clas-
504 sifier for a better certification performance while the others do not. Salman et al. and Lee use the
505 off-the-shelf classifier but without using the diffusion model. The most similar one compared with
506 us is Carlini et al., which also uses both the off-the-shelf diffusion model and classifier. The above
507 two settings mainly refer to Carlini et al. (2022), which makes us easier to compared with their
508 results.

509 E.3 Main Results for Certified Accuracy

510 We compare with Carlini et al. (2022) in a more fine-grained version. We provide results of certified
511 accuracy at different ϵ in Table A for CIFAR-10 and Table B for ImageNet. We include the accuracy
512 difference between ours and Carlini et al. (2022) in the bracket in Tables. We can observe from the
513 tables that the certified accuracy of our method outperforms Carlini et al. (2022) except $\epsilon = 0$ at
514 $\sigma = 0.25, 0.5$ for CIFAR-10.

Methods	Noise	Certified Accuracy at ϵ (%)					
		0.0	0.5	1.0	1.5	2.0	3.0
Carlini (Carlini et al., 2022)	$\sigma = 0.25$	77.0	71.0	0.0	0.0	0.0	0.0
	$\sigma = 0.5$	74.0	67.0	54.0	46.0	0.0	0.0
	$\sigma = 1.0$	59.0	53.0	49.0	38.0	29.0	22.0
Ours	$\sigma = 0.25$	80.0(+3.0)	76.0(+5.0)	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)
	$\sigma = 0.5$	75.0(+1.0)	72.0(+5.0)	62.0(+8.0)	49.0(+3.0)	0.0(+0.0)	0.0(+0.0)
	$\sigma = 1.0$	61.0(+2.0)	57.0(+4.0)	53.0(+4.0)	49.0(+11.0)	37.0(+8.0)	26.0(+4.0)

Table B: Certified accuracy compared with Carlini et al. (2022) for ImageNet at all σ . The numbers in the bracket are the difference of certified accuracy between two methods. Our diffusion model and classifier are the same as Carlini et al. (2022).

Datasets	Methods	Model	Certified Accuracy at ϵ (%)								
			0.0	0.25	0.5	0.75	Model	0.0	0.25	0.5	0.75
CIFAR-10	Carlini (Carlini et al., 2022)	ViT-B/16	93.0	76.0	57.0	47.0	WRN28-10	86.0	66.0	55.0	37.0
	Ours	ViT-B/16	92.0	82.0	69.0	56.0	WRN28-10	90.0	77.0	63.0	50.0
ImageNet	Carlini (Carlini et al., 2022)	BEiT	77.0	76.0	71.0	60.0	WRN50-2	73.0	67.0	57.0	48.0
	Ours	BEiT	80.0	78.0	76.0	71.0	WRN50-2	81.0	72.0	66.0	61.0

Table C: Certified accuracy of our method among different classifier. BEiT and ViT are pre-trained on a larger dataset ImageNet-22k and fine-tuned at ImageNet-1k and CIFAR-10 respectively. WideResNet is trained on ImageNet-1k for ImageNet and trained on CIFAR-10 from scratch for CIFAR-10.

515 E.4 Ablation study

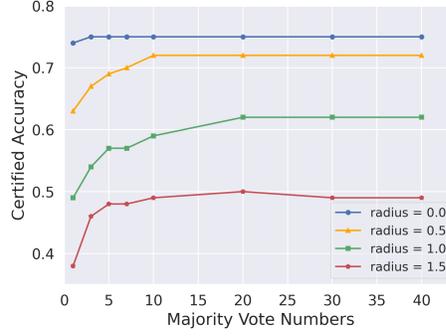
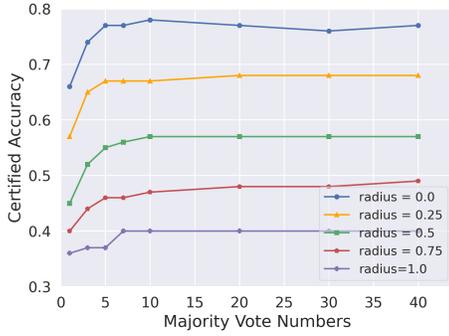
516 We conduct ablation study on different Voting samples. **Voting samples (K)** We first show how K
517 affects the certified accuracy. For efficiency, we select $b = 10$. We conduct experiments for both
518 datasets. We show the certified accuracy among different r at $\sigma = 0.25$ in Figure H. The results for
519 $\sigma = 0.5, 1.0$ and CIFAR-10 are shown in the Appendix E.5. Comparing with the baseline (Carlini
520 et al., 2022), we find that a larger majority vote number leads to a better certified accuracy. It verifies
521 that DensePure indeed benefits the adversarial robustness and making a good approximation of the
522 label with high density region requires a large number of voting samples. We find that our certified
523 accuracy will almost converge at $r = 40$. Thus, we set $r = 40$ for our experiments. The results with
524 other σ show the similar tendency.

525 **Fast sampling steps (b)** To investigate the role of b , we conduct additional experiments with $b \in$
526 $\{2, 5\}$ at $\sigma = 0.25$. The results on ImageNet are shown in Figure H and results for $\sigma = 0.5, 1.0$ and
527 CIFAR-10 are shown in the Appendix E.6. By observing results *with* majority vote, we find that a
528 larger b can lead to a better certified accuracy since a larger b generates images with higher quality.
529 By observing results *without* majority vote, the results show opposite conclusions where a larger b
530 leads to a lower certified accuracy, which contradicts to our intuition. We guess the potential reason
531 is that though more sampling steps can normally lead to better image recovery quality, it also brings
532 more randomness, increasing the probability that the reversed image locates into a data region with
533 the wrong label. These results further verify that majority vote is necessary for a better performance.

534 **Different architectures** One advantage of DensePure is to use the off-the-shelf classifier so that
535 it can plug in any classifier. We choose Convolutional neural network (CNN)-based architectures:
536 Wide-ResNet28-10 (Zagoruyko & Komodakis, 2016) for CIFAR-10 with 95.1% accuracy and Wide-
537 ResNet50-2 for ImageNet with 81.5% top-1 accuracy, at $\sigma = 0.25$. The results are shown in Table C
538 and Figure E in Appendix E.7. Results for more model architectures and σ of ImageNet are also
539 shown in Appendix E.7. We show that our method can enhance the certified robustness of any given
540 classifier trained on the original data distribution. Noticeably, although the performance of CNN-
541 based classifier is lower than Transformer-based classifier, DensePure with CNN-based model
542 as the classifier can outperform Carlini et al. (2022) with ViT-based model as the classifier (except
543 $\epsilon = 0$ for CIFAR-10).

544 E.5 Experiments for Voting Samples

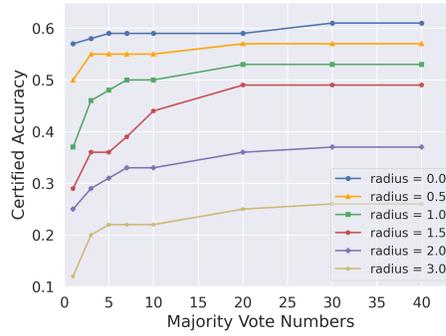
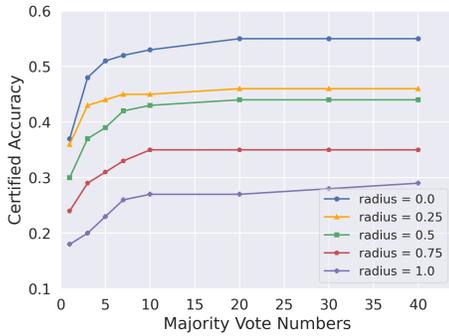
545 Here we provide more experiments with $\sigma \in \{0.5, 1.0\}$ and $b = 10$ for different voting samples K in
546 Figure A and Figure B. The results for CIFAR-10 is in Figure G. We can draw the same conclusion
547 mentioned in the main context .



CIFAR=10

ImageNet

Figure A: Certified accuracy among different vote numbers with different radius. Each line in the figure represents the certified accuracy among different vote numbers K with Gaussian noise $\sigma = 0.50$.



CIFAR=10

ImageNet

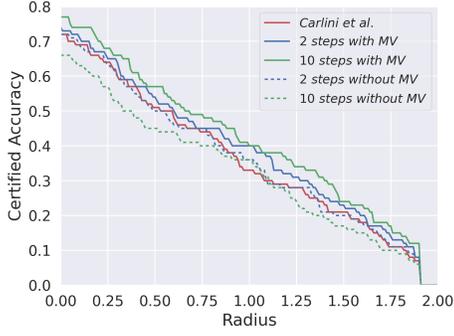
Figure B: Certified accuracy among different vote numbers with different radius. Each line in the figure represents the certified accuracy among different vote numbers K with Gaussian noise $\sigma = 1.00$.

548 **E.6 Experiments for Fast Sampling Steps**

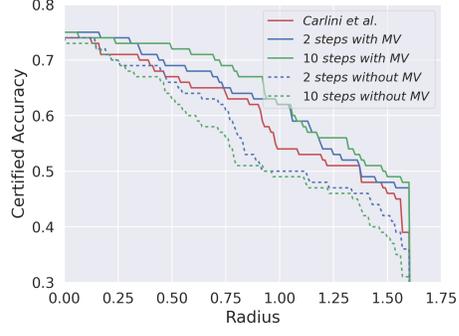
549 We also implement additional experiments with $b \in \{1, 2, 10\}$ at $\sigma = 0.5, 1.0$. The results are
 550 shown in Figure C and Figure D. The results for CIFAR-10 are in Figure G. We draw the same
 551 conclusion as mentioned in the main context.

552 **E.7 Experiments for Different Architectures**

553 We try different model architectures of ImageNet including Wide ResNet-50-2 and ResNet 152 with
 554 $b = 2$ and $K = 10$. The results are shown in Figure F. we find that our method outperforms (Carlini
 555 et al., 2022) for all σ among different classifiers.

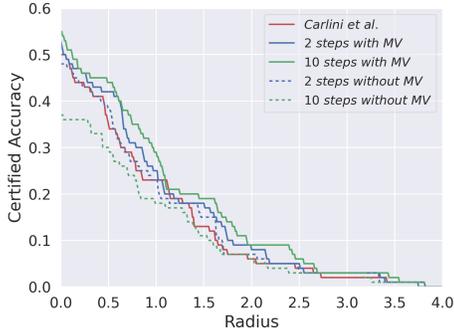


CIFAR=10

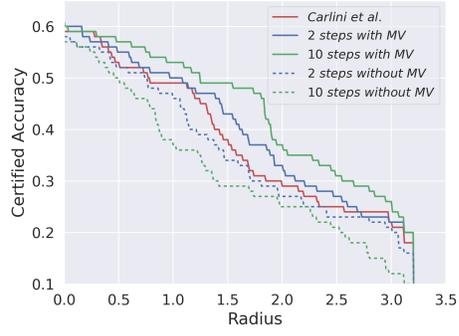


ImageNet

Figure C: Certified accuracy with different fast sampling steps b . Each line in the figure shows the certified accuracy among different L_2 adversarial perturbation bound with Gaussian noise $\sigma = 0.50$.

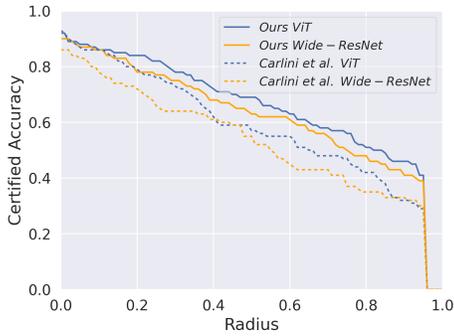


CIFAR=10

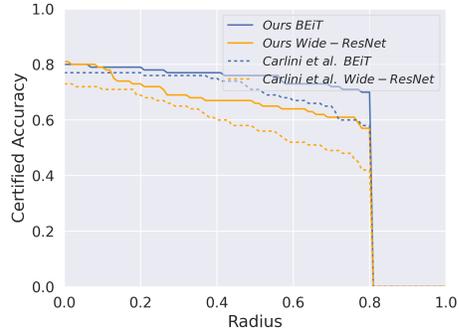


ImageNet

Figure D: Certified accuracy with different fast sampling steps b . Each line in the figure shows the certified accuracy among different L_2 adversarial perturbation bound with Gaussian noise $\sigma = 1.00$.

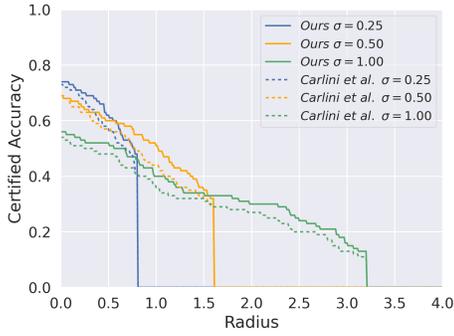


CIFAR=10

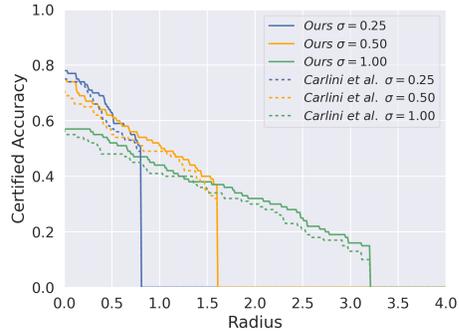


ImageNet

Figure E: Certified accuracy with different architectures. Each line in the figure shows the certified accuracy among different L_2 adversarial perturbation bound with Gaussian noise $\sigma = 0.25$.

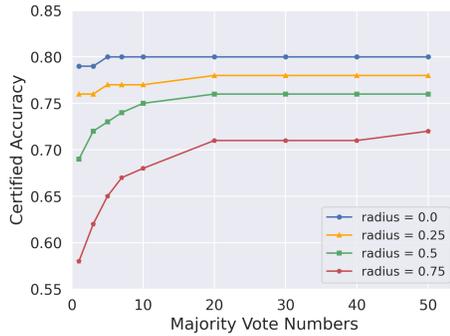


Wide ResNet-50-2

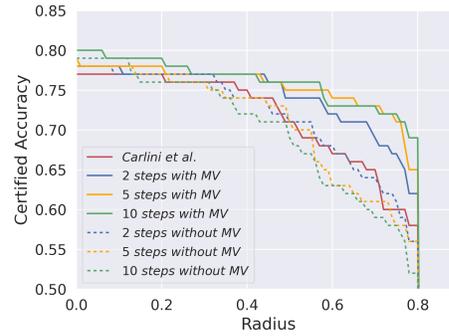


ResNet152

Figure F: Certified accuracy of ImageNet for different architectures. The lines represent the certified accuracy with different L_2 perturbation bound with different Gaussian noise $\sigma \in \{0.25, 0.50, 1.00\}$.



ImageNet



ImageNet

Figure G: Ablation study. The left image shows the certified accuracy among different vote numbers with different radius $\epsilon \in \{0.0, 0.25, 0.5, 0.75\}$. Each line in the figure represents the certified accuracy of our method among different vote numbers K with Gaussian noise $\sigma = 0.25$. The right image shows the certified accuracy with different fast sampling steps b . Each line in the figure shows the certified accuracy among different L_2 adversarial perturbation bound.

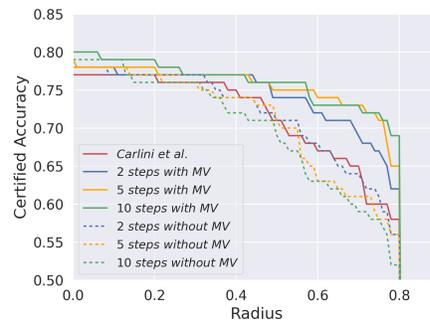
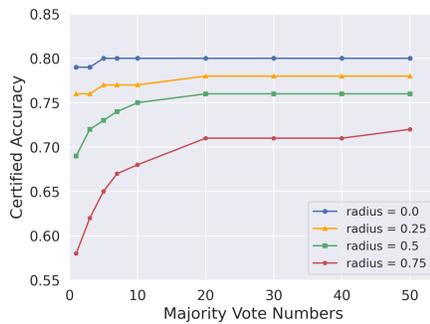


Figure H: Ablation study on ImageNet. The left image shows the certified accuracy among different vote numbers with different radius $\epsilon \in \{0.0, 0.25, 0.5, 0.75\}$. Each line in the figure represents the certified accuracy of our method among different vote numbers K with Gaussian noise $\sigma = 0.25$. The right image shows the certified accuracy with different fast sampling steps b . Each line in the figure shows the certified accuracy among different L_2 adversarial perturbation bound.