

## ETHICS STATEMENT

Our work can positively impact the society by improving the robustness and security of AI systems. We have not involved human subjects or data set releases; instead, we carefully follow the provided licenses of existing data and models for developing and evaluating our method.

## REPRODUCIBILITY STATEMENT

For theoretical analysis, all necessary assumptions are listed in B.1 and the complete proofs are included in B.2. The experimental setting and datasets are provided in section 5. The pseudo-code for DensePure is in C.1 and the fast sampling procedures are provided in C.2.

## REFERENCES

- Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- Nicholas Carlini, Florian Tramer, J Zico Kolter, et al. (certified!!) adversarial robustness for free! *arXiv preprint arXiv:2206.10550*, 2022.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1310–1320. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/cohen19c.html>.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *Advances in Neural Information Processing Systems*, 2019.
- Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *International Conference on Learning Representations*, 2020.
- Mitch Hill, Jonathan Craig Mitchell, and Song-Chun Zhu. Stochastic security: Adversarial defense using long-run dynamics of energy-based models. In *International Conference on Learning Representations*, 2021.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. URL <https://arxiv.org/abs/2006.11239>.
- Miklós Z Horváth, Mark Niklas Müller, Marc Fischer, and Martin Vechev. Boosting randomized smoothing with variance reduced classifiers. *arXiv preprint arXiv:2106.06946*, 2021.
- Jongheon Jeong and Jinwoo Shin. Consistency regularization for certified robustness of smoothed classifiers. *Advances in Neural Information Processing Systems*, 33:10558–10570, 2020.

- Jongheon Jeong, Sejun Park, Minkyu Kim, Heung-Chang Lee, Do-Guk Kim, and Jinwoo Shin. Smoothmix: Training confidence-calibrated smoothed classifiers for certified robustness. *Advances in Neural Information Processing Systems*, 34:30153–30168, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy (SP)*, pp. 656–672. IEEE, 2019.
- Kyungmin Lee. Provable defense by denoised smoothing with learned score function. In *ICLR Workshop on Security and Safety in Machine Learning Systems*, 2021.
- Matthew Mirman, Timon Gehr, and Martin Vechev. Differentiable abstract interpretation for provably robust neural networks. In *International Conference on Machine Learning*, pp. 3578–3586. PMLR, 2018.
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. Diffusion models for adversarial purification. In *International Conference on Machine Learning (ICML)*, 2022.
- Bernt Oksendal. *Stochastic differential equations: an introduction with applications*. Springer Science & Business Media, 2013.
- Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. In *International Conference on Learning Representations*, 2018a.
- Aditi Raghunathan, Jacob Steinhardt, and Percy S Liang. Semidefinite relaxations for certifying robustness to adversarial examples. In *NeurIPS*, 2018b.
- Hadi Salman, Jerry Li, Ilya Razenshteyn, Pengchuan Zhang, Huan Zhang, Sebastien Bubeck, and Greg Yang. Provably robust deep learning via adversarially trained smoothed classifiers. *Advances in Neural Information Processing Systems*, 32, 2019a.
- Hadi Salman, Greg Yang, Huan Zhang, Cho-Jui Hsieh, and Pengchuan Zhang. A convex relaxation barrier to tight robustness verification of neural networks. *Advances in Neural Information Processing Systems*, 32:9835–9846, 2019b.
- Hadi Salman, Mingjie Sun, Greg Yang, Ashish Kapoor, and J Zico Kolter. Denoised smoothing: A provable defense for pretrained classifiers. *Advances in Neural Information Processing Systems*, 33:21945–21957, 2020.
- Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-gan: Protecting classifiers against adversarial attacks using generative models. In *International Conference on Learning Representations*, 2018.
- Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *International Conference on Learning Representations*, 2018.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems*, 34:1415–1428, 2021a.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021b.
- Jiachen Sun, Weili Nie, Zhiding Yu, Z Morley Mao, and Chaowei Xiao. Pointdp: Diffusion-driven purification against adversarial attacks on 3d point cloud recognition. *arXiv preprint arXiv:2208.09801*, 2022.

- Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021.
- Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter. Beta-crown: Efficient bound propagation with per-neuron split constraints for neural network robustness verification. *Advances in Neural Information Processing Systems*, 34:29909–29921, 2021.
- Quanlin Wu, Hang Ye, and Yuntian Gu. Guided diffusion model for adversarial purification from random noise. *arXiv preprint arXiv:2206.10875*, 2022.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Runtian Zhai, Chen Dan, Di He, Huan Zhang, Boqing Gong, Pradeep Ravikumar, Cho-Jui Hsieh, and Liwei Wang. Macer: Attack-free and scalable robust training via maximizing certified radius. *arXiv preprint arXiv:2001.02378*, 2020.
- Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural network robustness certification with general activation functions. In *NeurIPS*, 2018.

## APPENDIX

Here is the appendix.

### 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
$\alpha_t$	scaling coefficient at time $t$
$\sigma_t^2$	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
$\epsilon_a$	adversarial attack
$\mathbf{x}_a$	adversarial sample
$\mathbf{x}_{a,t}$	scaled adversarial sample
$f(\cdot)$	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
$\mu_t, \nu_t$	path measure for $\{\hat{\mathbf{x}}_{\tau}\}_{\tau \in [0,t]}$ and $\{\mathbf{x}_{\tau}^{\theta}\}_{\tau \in [0,t]}$ respectively

## B MORE DETAILS ABOUT THEORETICAL ANALYSIS

### B.1 ASSUMPTIONS

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

### B.2 THEOREMS AND PROOFS

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

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

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

*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, which means

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

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

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

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

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 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,  $\mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$  is the robust region for data region  $\mathcal{G}(\mathbf{x}_0)$  under  $\mathcal{P}(\cdot; t)$  and  $f$ .

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

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  $\mathbf{x}_0$  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 should have

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

By Theorem 3.1, this is equivalent to

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

Re-organizing the above inequality, we obtain

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

Note that the order of  $\mathbf{x}_a$  is at most one in every term of the above inequality, so the inequality 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 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 obtain a convex  $\mathcal{D}_{\text{sub}}(\mathbf{x}_0; t)$ .

Then we prove part (ii).

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  $\mathbf{x}_a \in \mathcal{D}_{\text{sub}}(\tilde{\mathbf{x}}_0; t)$ . By part (i),  $\tilde{\mathbf{x}}_0$  has higher probability than all other points with different labels 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. Therefore,  $\mathcal{P}(\mathbf{x}_a; t)$  should have the same label as  $\mathbf{x}_0$ . On the other hand, if  $\mathbf{x}_a \notin \mathcal{D}(\mathcal{G}(\mathbf{x}_0); t)$ , 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  $\mathbf{x}_0$ ,  $\mathbb{P}(\hat{\mathbf{x}}_0 = \tilde{\mathbf{x}}_1 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t}) > \mathbb{P}(\hat{\mathbf{x}}_0 = \tilde{\mathbf{x}}_0 | \hat{\mathbf{x}}_t = \mathbf{x}_{a,t})$ . In other words,  $\mathcal{P}(\mathbf{x}_a; t)$  would have different label from  $\mathbf{x}_0$ .  $\square$

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

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

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,

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

$\mathbf{s}_\theta(\mathbf{x}, \tau)$  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 used in the training score-based diffusion models.

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

(2013):

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

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

## C MORE DETAILS ABOUT DENSEPURE

### C.1 PSEUDO-CODE

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  $\mathbf{x}_a = \mathbf{x}_0 + \boldsymbol{\epsilon}_a$
  - 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  $\sigma$
  - 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  $\mathbf{x}_{\text{rs}} = \mathbf{x}_0 + \boldsymbol{\epsilon}$ ,  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
  - 4: Choose schedule  $S^b$ , get  $\hat{\mathbf{x}}_0^i \leftarrow \text{rev}(\sqrt{\bar{\alpha}_n} \mathbf{x}_{\text{rs}})_i$ ,  $i = 1, 2, \dots, K$  with Fast Sampling
  - 5:  $\hat{y} = \mathbf{MV}(\{f(\hat{\mathbf{x}}_0^1), \dots, f(\hat{\mathbf{x}}_0^K)\}) = \arg \max_c \sum_{i=1}^K \mathbf{1}\{f(\hat{\mathbf{x}}_0^i) = c\}$
- 

### C.2 DETAILS ABOUT FAST SAMPLING

Applying single-step operation  $n$  times is a time-consuming process. In order to reduce the time complexity, we follow the method used in (Nichol & Dhariwal, 2021) and sample a subsequence  $S^b$  with  $b$  values (i.e.,  $S^b = \{n, \underbrace{\lfloor n - \frac{n}{b} \rfloor, \dots, 1}_b\}$ , where  $S_j^b$  is the  $j$ -th element in  $S^b$  and  $S_j^b =$

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

$S_j = j$  is the  $j$ -th element in  $S$ ).

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-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.,  $\bar{\alpha}_i^{S^b} = \bar{\alpha}_{S_i^b} = \bar{\alpha}_{\lfloor n - \frac{jn}{b} \rfloor}$  is the  $i$ -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  $\tilde{\beta}^{S^b} = \{\tilde{\beta}_1^{S^b}, \tilde{\beta}_2^{S^b}, \dots, \tilde{\beta}_i^{S^b}, \dots, \tilde{\beta}_b^{S^b}\}$  schedules, where  $\beta_{S_i^b} = \beta_i^{S^b} = 1 - \frac{\bar{\alpha}_i^{S^b}}{\bar{\alpha}_{i-1}^{S^b}}$ ,  $\tilde{\beta}_{S_i^b} = \tilde{\beta}_i^{S^b} = \frac{1 - \bar{\alpha}_{i-1}^{S^b}}{1 - \bar{\alpha}_i^{S^b}} \beta_{S_i^b}$ . With these new schedules, we can use  $b$  times reverse steps to calculate

Methods	Noise	Certified Accuracy at $\epsilon$ (%)				
		0.0	0.25	0.5	0.75	1.0
Carlini (Carlini et al., 2022)	$\sigma = 0.25$	<b>88.0</b>	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
<b>Ours</b>	$\sigma = 0.25$	87.6(-0.4)	<b>76.6(+2.8)</b>	<b>64.6(+8.4)</b>	<b>50.4(+8.8)</b>	0.0(+0.0)
	$\sigma = 0.5$	73.6(-0.6)	65.4(+3.4)	55.6(+5.2)	46.0(+5.8)	<b>37.4(+6.4)</b>
	$\sigma = 1.0$	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  $\sigma$ . 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).

$\hat{x}_0 = \underbrace{\text{Reverse}(\cdots \text{Reverse}(\text{Reverse}(x_n; S_b^b); S_{b-1}^b); \cdots; 1)}_b$ . Since  $\Sigma_\theta(x_{S_i^b}, S_i^b)$  is parameterized as a range between  $\beta^{S^b}$  and  $\tilde{\beta}^{S^b}$ , it will automatically be rescaled. Thus,  $\hat{x}_{S_{i-1}^b} = \text{Reverse}(\hat{x}_{S_i^b}; S_i^b)$  is equivalent to sample  $x_{S_{i-1}^b}$  from  $\mathcal{N}(x_{S_{i-1}^b}; \mu_\theta(x_{S_i^b}, S_i^b), \Sigma_\theta(x_{S_i^b}, S_i^b))$ .

## D MORE EXPERIMENTAL DETAILS AND RESULTS

### D.1 IMPLEMENTATION DETAILS

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

### D.2 BASELINES.

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

### D.3 MAIN RESULTS FOR CERTIFIED ACCURACY

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

### D.4 EXPERIMENTS FOR VOTING SAMPLES

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



Methods	Noise	Certified Accuracy at $\epsilon(\%)$					
		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
<b>Ours</b>	$\sigma = 0.25$	<b>80.0(+3.0)</b>	<b>76.0(+5.0)</b>	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)	<b>62.0(+8.0)</b>	<b>49.0(+3.0)</b>	0.0(+0.0)	0.0(+0.0)
	$\sigma = 1.0$	61.0(+2.0)	57.0(+4.0)	53.0(+4.0)	<b>49.0(+11.0)</b>	<b>37.0(+8.0)</b>	<b>26.0(+4.0)</b>

Table B: Certified accuracy compared with Carlini et al. (2022) for ImageNet at all  $\sigma$ . 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).

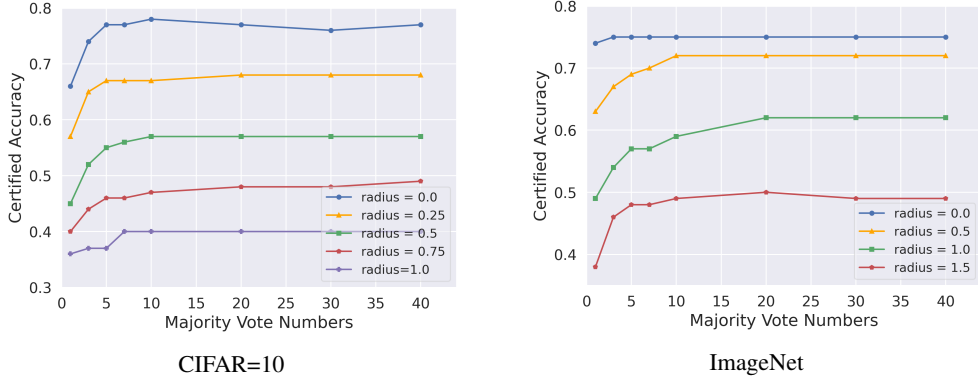


Figure A: Certified accuracy among different vote numbers with different radius. Each line in the figure represents the certified accuracy among different vote numbers  $K$  with Gaussian noise  $\sigma = 0.50$ .

#### D.5 EXPERIMENTS FOR FAST SAMPLING STEPS

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

#### D.6 EXPERIMENTS FOR DIFFERENT ARCHITECTURES

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

#### D.7 EXPERIMENTS FOR RANDOMIZED SMOOTHING WITHOUT DIFFUSION MODEL

To explore randomized smoothing without diffusion model, we directly remove the diffusion model from our pipeline and conduct additional experiments.

First, we remove the diffusion model and perform randomized smoothing only on the pretrained classifier we used in DensePure (i.e., ViT-B/16 for CIFAR-10 and BEiT for ImageNet). The results are shown in Table C and Table D. The number in the bracket is calculated by the robust accuracy of pretrained classifier - the robust accuracy of DensePure. We can conclude from the table that without the help of diffusion models, neither ViT nor BEiT could reach high certified accuracy.

Second, we conduct additional experiments to fairly compare with randomized smoothing without diffusion models under majority vote settings. Specifically, we activate droppath in BEiT at the inference stage to support majority votes. The other settings are the same as DensePure. The results are shown in Table E. The number in the bracket is calculated by the robust accuracy of BEiT with majority votes - the robust accuracy of DensePure. We find that simply performing majority votes on the BEiT classifier will not result in higher certified robustness.

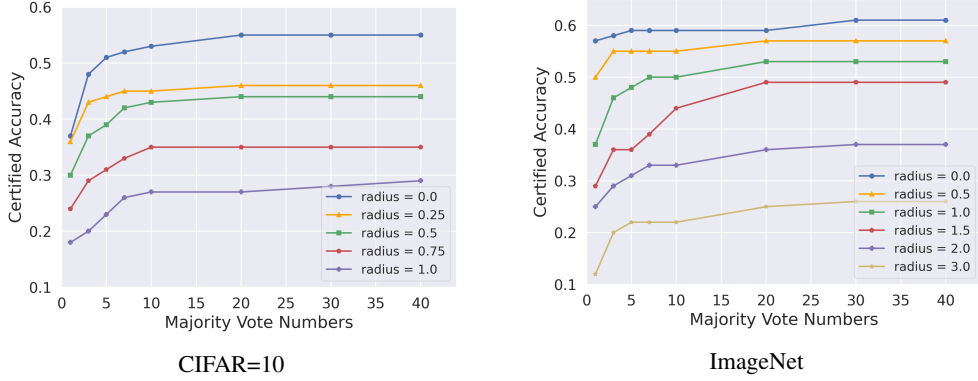


Figure B: Certified accuracy among different vote numbers with different radius. Each line in the figure represents the certified accuracy among different vote numbers  $K$  with Gaussian noise  $\sigma = 1.00$ .

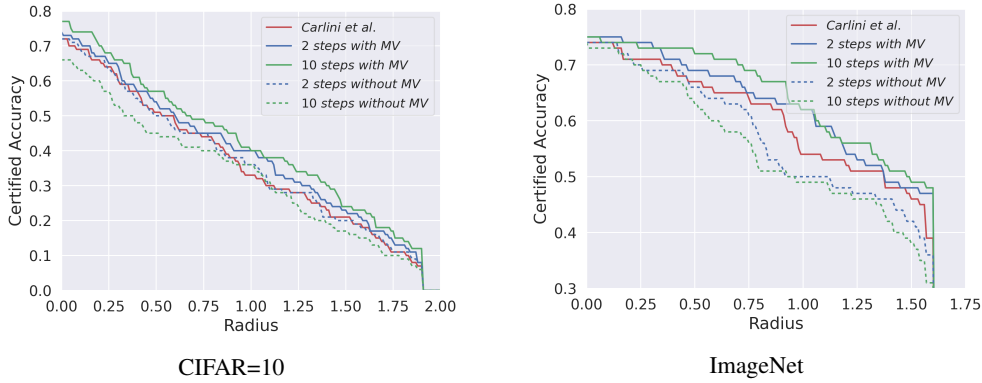


Figure C: Certified accuracy with different fast sampling steps  $b$ . Each line in the figure shows the certified accuracy among different  $L_2$  adversarial perturbation bound with Gaussian noise  $\sigma = 0.50$ .

Third, to compare with randomized smoothing without diffusion model, we also evaluate certified accuracy with Gaussian augmentation-trained ViT models on CIFAR-10. The results shown in the table F prove that DensePure can still achieve higher certified accuracy than randomized smoothing on even Gaussian augmented models without diffusion models. The numbers in the bracket are the difference between the robust accuracy of Gaussian augmentation randomized smoothing and DensePure.

#### D.8 EXPERIMENTS FOR K-CONSENSUS AGGREGATION

In K-Consensus Aggregation, if the classification results of the  $K$  consecutive reversed samples are the same, an early stop will be triggered. Here We calculate certified robustness for 100 subsamples of CIFAR-10 and ImageNet with 2 sampling steps, a maximum 10 majority votes and consensus threshold  $k=3$ . Results are shown in Table G and Table H. The column of "Avg MV" in the tables means the average of the actual number of majority votes required for our algorithm. For instance, if the predicted labels of the first 3 reversed samples are the same, the actual majority vote numbers will be 3. The numbers in the bracket are the difference between certified accuracy w/o K-Consensus Aggregation.

#### D.9 EXPERIMENTS FOR CERTIFIED ACCURACY WITH LESS SAMPLING STEPS AND VOTE NUMBERS

We also conduct additional experiments with 2 sampling steps and 5 majority votes. The results are shown in Table I. We find that our method still achieves better results than the existing method.

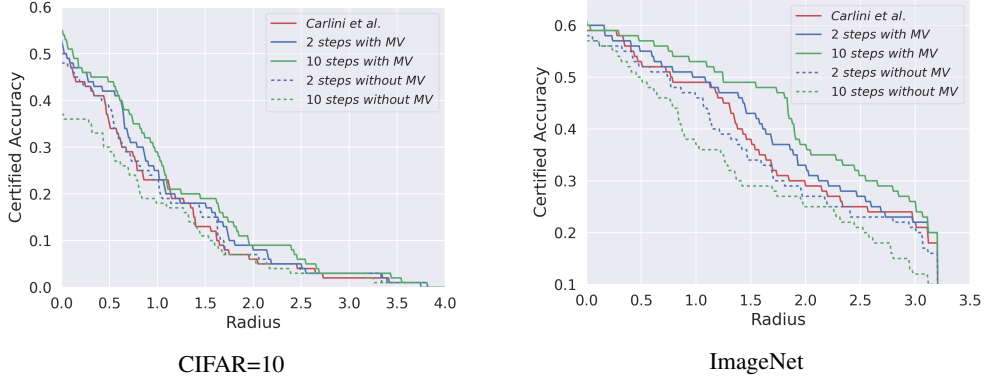


Figure D: Certified accuracy with different fast sampling steps  $b$ . Each line in the figure shows the certified accuracy among different  $L_2$  adversarial perturbation bound with Gaussian noise  $\sigma = 1.00$ .

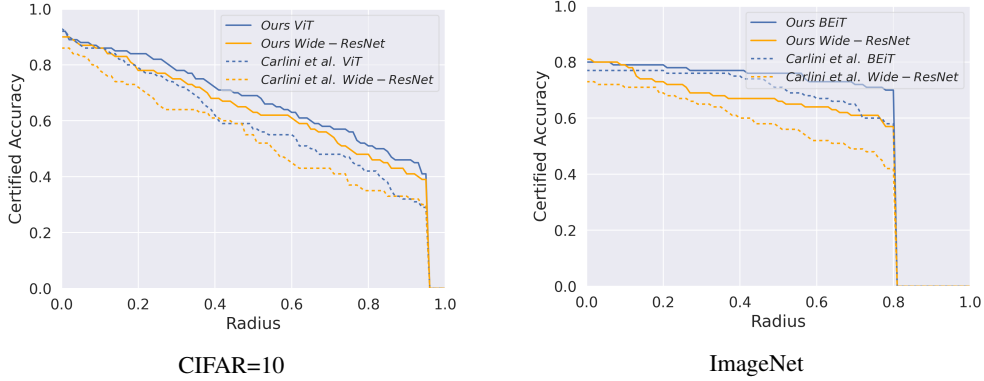


Figure E: Certified accuracy with different architectures. Each line in the figure shows the certified accuracy among different  $L_2$  adversarial perturbation bound with Gaussian noise  $\sigma = 0.25$ .

#### D.10 EXPERIMENTS FOR DENSEPURE 500 TEST SAMPLING NUMBER RESULTS ON IMAGENET

We increase the ImageNet test sampling number from 100 to 500 and update the experiment results in Table J and Table K.

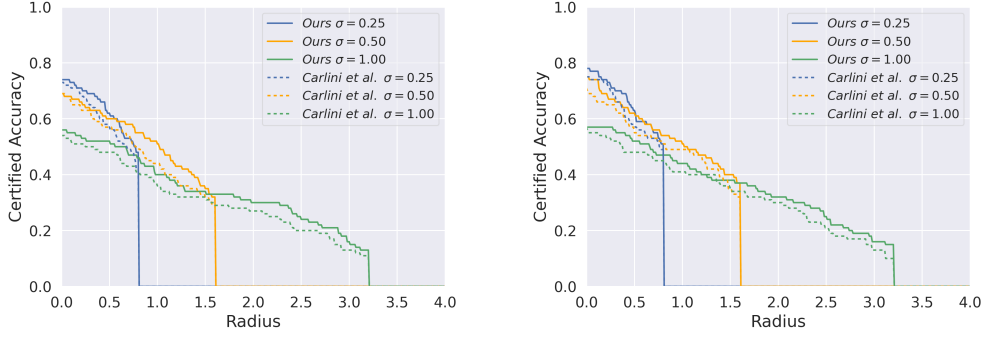


Figure F: Certified accuracy of ImageNet for different architectures. The lines represent the certified accuracy with different  $L_2$  perturbation bound with different Gaussian noise  $\sigma \in \{0.25, 0.50, 1.00\}$ .

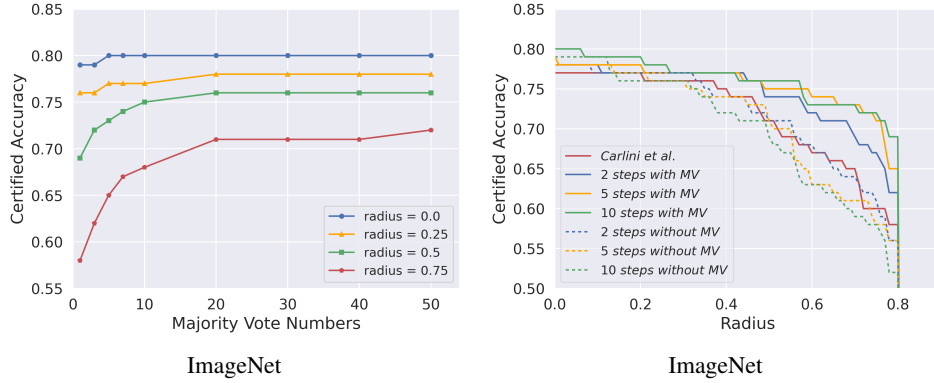


Figure G: Ablation study. The left image shows the certified accuracy among different vote numbers with different radius  $\epsilon \in \{0.0, 0.25, 0.5, 0.75\}$ . Each line in the figure represents the certified accuracy of our method among different vote numbers  $K$  with Gaussian noise  $\sigma =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

Noise	Certified Accuracy at $\epsilon(\%)$				
	0.0	0.25	0.5	0.75	1.0
$\sigma = 0.25$	20.8(-66.8)	7.4(-69.2)	1.8(-62.8)	0.2(-50.2)	0.0(+0.0)
$\sigma = 0.5$	11.6(-62.0)	6.6(-58.8)	3.8(-51.8)	1.2(-44.8)	0.2(-37.2)
$\sigma = 1.0$	10.6(-44.4)	10.6(-37.4)	9.4(-31.4)	9.4(-23.6)	9.4(-18.8)

Table C: Certified accuracy of randomized smoothing on pretrained classifier ViT-B/16 at all  $\sigma$  for CIFAR-10

Noise	Certified Accuracy at $\epsilon(\%)$					
	0.0	0.5	1.0	1.5	2.0	3.0
$\sigma = 0.25$	73.2(-10.8)	55.8(-22.0)	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)
$\sigma = 0.5$	7.8(-72.4)	4.6(-71.0)	3.2(-63.8)	1.0(-53.6)	0.0(+0.0)	0.0(+0.0)
$\sigma = 1.0$	0.0(-67.8)	0.0(-61.4)	0.0(-55.6)	0.0(-50.0)	0.0(-42.2)	0.0(-25.8)

Table D: Certified accuracy of randomized smoothing on pretrained classifier BEiT at all  $\sigma$  for ImageNet

Noise	Certified Accuracy at $\epsilon$ (%)					
	0.0	0.5	1.0	1.5	2.0	3.0
$\sigma = 0.25$	73.8(-10.2)	58.0(-19.8)	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)
$\sigma = 0.5$	9.0(-71.2)	7.0(-68.6)	4.0(-63.0)	2.0(-52.6)	0.0(+0.0)	0.0(+0.0)
$\sigma = 1.0$	0.0(-67.8)	0.0(-61.4)	0.0(-55.6)	0.0(-50.0)	0.0(-42.2)	0.0(-25.8)

Table E: Certified accuracy of randomized smoothing on droppatch activated BEiT with 10 majority votes at all  $\sigma$  for ImageNet

Noise	Certified Accuracy at $\epsilon$ (%)				
	0.0	0.25	0.5	0.75	1.0
$\sigma = 0.25$	88.2(+0.6)	71.4(-5.2)	53.2(-11.4)	35.2(-15.2)	0.0(+0.0)
$\sigma = 0.5$	69.8(-3.8)	60.0(-5.4)	48.4(-7.2)	37.2(-8.8)	27.2(-10.2)
$\sigma = 1.0$	49.0(-6.0)	41.8(-6.0)	34.0(-6.8)	27.0(-6.0)	22.0(-6.2)

Table F: Certified accuracy of randomized smoothing on Gaussian augmentation-trained ViT at all  $\sigma$  on CIFAR-10

Noise	Certified Accuracy at $\epsilon$ (%)					Avg MV
	0.0	0.25	0.5	0.75	1.0	
$\sigma = 0.25$	92(+0.0)	77(+0.0)	60(+0.0)	48(-1.0)	0(+0.0)	3.84
$\sigma = 0.5$	74(+0.0)	65(+0.0)	53(-1.0)	45(+0.0)	40(+0.0)	4.43
$\sigma = 1.0$	53(+0.0)	46(+0.0)	42(+0.0)	31(+0.0)	25(+0.0)	5.49

Table G: Certified accuracy and average majority votes with 2 sample steps and  $k = 3$  consensus threshold at all  $\sigma$  for CIFAR-10.

Noise	Certified Accuracy at $\epsilon$ (%)						Avg MV
	0.0	0.5	1.0	1.5	2.0	3.0	
$\sigma = 0.25$	78(+0.0)	74(+0.0)	0(+0.0)	0(+0.0)	0(+0.0)	0(+0.0)	3.34
$\sigma = 0.5$	75(+0.0)	69(+0.0)	61(+0.0)	47(+0.0)	0(+0.0)	0(+0.0)	3.89
$\sigma = 1.0$	60(+0.0)	54(+0.0)	50(+0.0)	41(+0.0)	32(+0.0)	23(+0.0)	5.23

Table H: Certified accuracy and average majority votes with 2 sample steps and  $k = 3$  consensus threshold at all  $\sigma$  for ImageNet.

Noise	Certified Accuracy at $\epsilon$ (%)										
	CIFAR-10					ImageNet					
	0.0	0.25	0.5	0.75	1.0	0.0	0.5	1.0	1.5	2.0	3.0
$\sigma = 0.25$	87.6	74.8	59.2	44.6	0.0	78	74	0	0	0	0
$\sigma = 0.50$	73.2	62.6	52.6	41.8	34.0	75	69	58	47	0	0
$\sigma = 1.00$	53.4	44.0	35.8	30.2	24.4	60	54	49	39	30	22

Table I: Certified accuracy with 2 sampling steps and 5 vote numbers at all  $\sigma$  for both CIFAR-10 and ImageNet

Method	Off-the-shelf	Certified Accuracy at $\epsilon$ (%)									
		CIFAR-10					ImageNet				
		0.25	0.5	0.75	1.0	1.0	0.5	1.0	1.5	2.0	3.0
PixelDP (Lecuyer et al., 2019)	$\times$	(71.0)22.0	(44.0)2.0	-	-	(33.0)16.0	-	-	-	-	-
RS (Cohen et al., 2019)	$\times$	(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)	$\times$	(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 & Shin, 2020)	$\times$	(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)	$\times$	(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)	$\times$	(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	(44.6)28.6	(38.6)21.2	-
SmoothMix (Jeong et al., 2021)	$\checkmark$	(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)	$\checkmark$	(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)	$\checkmark$	60.0	42.0	28.0	19.0	41.0	24.0	11.0	-	-	-
Carlini (Carlini et al., 2022)	$\checkmark$	(88.0)73.8	(88.0)56.2	(88.0)41.6	(74.2)31.0	(82.0)74.0	(77.2.0)59.8	(77.2)47.0	(64.6)31.0	(64.6)19.0	-
<b>Ours</b>	$\checkmark$	(87.6) <b>76.6</b>	(87.6) <b>64.6</b>	(87.6)50.4	(73.6)37.4	(84.0) <b>77.8</b>	(80.2) <b>67.0</b>	(80.2) <b>54.6</b>	(67.8) <b>42.2</b>	(67.8) <b>25.8</b>	-

Table J: 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).

Methods	Noise	Certified Accuracy at $\epsilon$ (%)					
		0.0	0.5	1.0	1.5	2.0	3.0
Carlini (Carlini et al., 2022)	$\sigma = 0.25$	82.0	74.0	0.0	0.0	0.0	0.0
	$\sigma = 0.5$	77.2	71.8	59.8	47.0	0.0	0.0
	$\sigma = 1.0$	64.6	57.8	49.2	40.6	31.0	19.0
<b>Ours</b>	$\sigma = 0.25$	<b>84.0(+2.0)</b>	<b>77.8(+3.8)</b>	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)	0.0(+0.0)
	$\sigma = 0.5$	80.2(+3.0)	75.6(+3.8)	<b>67.0(+7.2)</b>	<b>54.6(+7.6)</b>	0.0(+0.0)	0.0(+0.0)
	$\sigma = 1.0$	67.8(+3.2)	61.4(+3.6)	55.6(+6.4)	50.0(+9.4)	<b>42.2(+11.2)</b>	<b>25.8(+6.8)</b>

Table K: Certified accuracy compared with Carlini et al. (2022) for ImageNet at all  $\sigma$ . 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).