
Understanding and Enhancing the Robustness of Diffusion-Based Purification

Anonymous Author(s)

Affiliation
Address
email

Abstract

Diffusion-Based Purification (DBP) has emerged as an effective defense mechanism against adversarial attacks. Traditionally, the efficacy of DBP has been attributed to the forward diffusion process, which narrows the distribution gap between clean and adversarial images through the addition of Gaussian noise. Although theoretical studies support this explanation, to what extent it contributes to robustness remains unclear. In this paper, we argue that the inherent stochasticity in the DBP process is the primary driver of its robustness. To explore this, we introduce a novel Deterministic White-Box (DW-box) evaluation framework to assess robustness and analyze the attack trajectories and loss landscapes. Our findings suggest that DBP models primarily leverage stochasticity to evade effective attack directions, rather than directly neutralizing adversarial perturbations. To further enhance DBP robustness, we integrate adversarial perturbations into diffusion training, propose Rank-Based Gaussian Mapping (RBGM) to make perturbations more compatible with the diffusion models, and introduce Adversarial Denoising Diffusion Training (ADDT) to strengthen diffusion models with classifier-guided perturbations. Empirical evidence demonstrates the effectiveness of ADDT.

1 Introduction

Deep learning has achieved remarkable success in various domains, including computer vision [He et al., 2016], natural language processing [OpenAI, 2023], and speech recognition [Radford et al., 2022]. However, in this flourishing landscape, the persistent specter of adversarial attacks casts a shadow over the reliability of these neural models. Adversarial attacks for a vision model involve injecting imperceptible perturbations into input images to trick models into producing false outputs with high confidence [Goodfellow et al., 2015; Szegedy et al., 2014]. This inspires a large amount of research on adversarial defense [Zhang et al., 2019; Samangouei et al., 2018; Shafahi et al., 2019a; Wang et al., 2023].

Diffusion-based purification (DBP) [Nie et al., 2022] has recently gained recognition as a powerful defense mechanism against a range of adversarial attacks, exploiting the capabilities of available diffusion models. The conventional view suggests that the robustness provided by DBP is primarily due to the forward diffusion process, which narrows the distribution gap between clean and adversarial

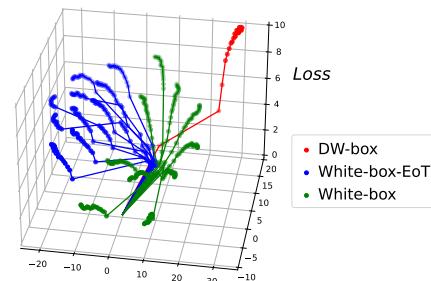


Figure 1: Comparison of attack trajectories under different evaluation settings. The visualization shows the effectiveness of attacks under the Deterministic White-box (DW-box) setting, while the attack trajectory under the standard White-box setting is less effective and deviates significantly from the DW-box trajectory.

38 images through the application of Gaussian noise [Wang et al., 2022; Nie et al., 2022]. However,
39 although the reduction of the distribution gap is theoretically proven, its actual contribution to im-
40 proving robustness has not been sufficiently investigated. Thus, it remains an open question of which
41 factor contributes to DBP robustness most.

42 In light of this, we introduce an alternative perspective that highlights the role of stochasticity
43 throughout the DBP process as a key contributor to its robustness, challenging the traditional focus
44 on the forward diffusion process. To evaluate the impact of stochasticity, we employ a Deterministic
45 White-box (DW-box) attack setting where the attacker has complete knowledge of both the model
46 parameters and the stochastic elements. Our findings reveal that DBP models significantly lose their
47 robustness when the process is entirely deterministic to the attacker, thereby emphasizing the critical
48 importance of stochasticity. Further investigations into attack trajectories and the loss landscape
49 demonstrate that DBP models do not inherently counter adversarial perturbations as effectively as
50 models trained with adversarial training (AT). Instead, they rely on stochasticity to circumvent the
51 effective attack direction. This dependency is discussed in detail in Section 4.2 and visually depicted
52 in Figure 1.

53 To enhance the robustness of DBP models, we introduce Adversarial Denoising Diffusion Train-
54 ing (ADDT). ADDT employs an iterative two-step approach: the Classifier-Guided Perturbation
55 Optimization (CGPO) step generates adversarial perturbations, while the training step updates the
56 diffusion model parameters using these perturbations. To align the perturbations more closely with
57 the diffusion framework, we introduce Rank-Based Gaussian Mapping (RBGM), which adapts the
58 adversarial perturbations to be more Gaussian-like.

59 Our main contributions are as follows:

- 60 • We present a novel perspective on DBP robustness, emphasizing the critical role of stochas-
61 ticity and challenging the conventional belief that robustness primarily stems from reducing
62 the distribution gap via the forward diffusion process.
- 63 • We introduce a new deterministic white-box attack setting and show that DBP models
64 depend on stochastic elements to avoid effective attack directions, demonstrating distinct
65 properties compared to models trained with Adversarial Training (AT).
- 66 • We develop Adversarial Denoising Diffusion Training (ADDT), which enhances DBP
67 models' robustness, and introduce Rank-Based Gaussian Mapping (RBGM) to adapt pertur-
68 bations to be more Gaussian-like, aligning them with the diffusion framework. Empirical
69 validation confirms that ADDT achieves a robust accuracy improvement of up to 6% com-
70 pared to conventional DBP models.

71 2 Related Work

72 **Adversarial Training.** First introduced by Madry et al. [2018], adversarial training (AT) seeks to
73 develop a robust classifier by incorporating adversarial examples into the training process. It has
74 nearly become the de facto standard for enhancing the adversarial robustness of neural networks
75 [Gowal et al., 2020; Rebuffi et al., 2021; Athalye et al., 2018]. Recent advances in AT harness the
76 generative power of diffusion models to augment training datasets and prevent AT from overfitting
77 [Gowal et al., 2021; Wang et al., 2023]. However, the application of AT to DBP methods has not
78 been thoroughly explored.

79 **Adversarial Purification.** Adversarial purification utilizes generative models to remove adversarial
80 perturbation from inputs before they are processed by downstream models. Traditionally, generative
81 adversarial networks (GANs) [Samangouei et al., 2018] or autoregressive models [Song et al., 2018]
82 are employed as the purifier model. More recently, diffusion models have been introduced for
83 adversarial purification, in a technique termed diffusion-based purification (DBP), and have shown
84 promising results [Song and Ermon, 2019; Ho et al., 2020; Song et al., 2020a; Nie et al., 2022;
85 Wang et al., 2022; Wu et al., 2022; Xiao et al., 2022]. The robustness of DBP models is often
86 attributed to the wash-out effect of Gaussian noise introduced during the forward diffusion process.
87 Nie et al. [2022] propose that the forward process results in a reduction of the Kullback-Leibler (KL)
88 divergence between the distributions of clean and adversarial images. Gao et al. [2022] suggest that
89 while the forward diffusion process improves robustness by reducing model invariance, the backward
90 process restores this invariance, thereby undermining robustness. However, these theories explaining
91 the robustness of DBP models lack substantial experimental support.

92 **3 Preliminaries**

93 **Adversarial Training.** Adversarial training aims to create a robust model by including adversarial
 94 samples during training [Madry et al., 2018]. This approach can be formulated as a min-max problem,
 95 where we first generate adversarial samples (the maximization) and then adjust the parameters to
 96 resist these adversarial samples (the minimization). Mathematically, this is represented as:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\max_{\delta \in B} L(f(\boldsymbol{\theta}, \mathbf{x} + \delta), y)], \quad (1)$$

97 where L is the loss function, f is the classifier, $(\mathbf{x}, y) \sim \mathcal{D}$ denotes sampling training data from
 98 distribution \mathcal{D} , and B defines the set of permissible perturbation δ .

99 **Diffusion Models.** Denoising Diffusion Probabilistic Models (DDPM) [Ho et al., 2020] and De-
 100 noising Diffusion Implicit Models (DDIM) [Song et al., 2020a] simulate a gradual transformation in
 101 which noise is added to images and then removed to restore the original image. The forward process
 102 can be represented as:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}), \quad (2)$$

103 where \mathbf{x}_0 is the original image and \mathbf{x}_t is the noisy image. $\bar{\alpha}_t$ is the cumulative noise level at step t
 104 ($1 < t \leq T$, where T is the number of diffusion training steps). The model optimizes the parameters
 105 $\boldsymbol{\theta}$ by minimizing the distance between the actual and predicted noise:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x}_0, t, \epsilon} [\|\epsilon - \epsilon_{\boldsymbol{\theta}}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|_2^2], \quad (3)$$

106 where $\epsilon_{\boldsymbol{\theta}}$ is the model's noise prediction, with $\epsilon_{\boldsymbol{\theta}}$, we can predict $\hat{\mathbf{x}}_0$ in a single step:

$$\hat{\mathbf{x}}_0 = (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\boldsymbol{\theta}^*}(\mathbf{x}_t, t)) / \sqrt{\bar{\alpha}_t}, \quad (4)$$

107 where $\hat{\mathbf{x}}_0$ is the recovered image. DDPM typically takes an iterative approach to restore the image,
 108 removing a small amount of Gaussian noise at a time:

$$\hat{\mathbf{x}}_{t-1} = \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\boldsymbol{\theta}^*}(\mathbf{x}_t, t) \right) / \sqrt{1 - \beta_t} + \sqrt{\beta_t} \epsilon, \quad (5)$$

109 where β_t is the noise level at step t , $\hat{\mathbf{x}}_{t-1}$ is the recovered image in step $t - 1$, ϵ is sampled from
 110 $\mathcal{N}(0, \mathbf{I})$. DDIM proposes to speed up the denoising process by skipping certain intermediate steps.
 111 Recent work suggests that DDPM may also benefit from a similar approach [Nichol and Dhariwal,
 112 2021]. Score SDEs [Song et al., 2020b] give a score function view of DDPM and further lead to the
 113 derivations of DDPM++ (VPSDE) and EDM [Karras et al., 2022].

114 **Diffusion-Based Purification (DBP).** DBP uses diffusion models to remove adversarial perturbation
 115 from images. Instead of using a complete diffusion process between the clean image and pure
 116 Gaussian noise (between $t = 0$ and $t = T$), they first diffuse \mathbf{x}_0 to a predefined timestep $t = t^*$ ($t^* <$
 117 T) via Equation (2), and recover the image $\hat{\mathbf{x}}_0$ via the reverse diffusion process in Equation (5).

118 **4 Stochasticity-Driven Robustness**

119 **4.1 Stochasticity as the Main Factor of DBP Robustness**

120 Traditional perspectives, as discussed in Section 2, primarily attribute the robustness of DBP models
 121 to the forward diffusion process, which introduces Gaussian noise into both clean and adversarial
 122 images, thereby reducing the difference between their distributions [Wang et al., 2022; Nie et al.,
 123 2022]. Although supported by theoretical studies, to what extent this contributes to the practical
 124 robustness remains an open question.

125 Considering the stochastic nature of DBP, we provide an alternative perspective that the robustness of
 126 DBP models may not only derive from the forward diffusion process but also depend on the stochastic
 127 elements integrated throughout the DBP process. To test this hypothesis, we evaluate the contributions
 128 of both the *forward diffusion process* and the *stochasticity throughout the processes* to the robustness
 129 of several DBP models. In particular, we focus on two models: DDPM, which incorporates Gaussian
 130 noise in both the forward and backward processes, and DDIM, which incorporates Gaussian noise in
 131 the forward process only. We adopt DDPM/DDIM within the DiffPure framework [Nie et al., 2022]
 132 and denote the implementations by DP_{DDPM} and DP_{DDIM}, respectively.

133 To isolate the role of stochasticity, we introduce a new attack scenario: the **Deterministic White-Box**
 134 (DW-box) setting. This contrasts with the traditional White-box setting, where the attacker only has
 135 knowledge of the model parameters and must rely on random sampling for the stochastic elements
 136 involved in the evaluation. In the DW-box setting, the attacker has full knowledge of both the model
 137 parameters and the specific states of stochastic elements used in the evaluation. This setting makes the
 138 DBP process deterministic to the attacker. We further differentiate the settings based on the attacker’s
 139 knowledge of the stochastic elements: no knowledge (White-box), full knowledge in the forward
 140 process (DW_{Fwd}-box), and full knowledge in both the forward and reverse processes (DW_{Both}-box).
 141 We detail the differences between these settings in Appendix B.

142 According to the conventional understanding, which at-
 143 tributes the robustness of DBP models solely to the for-
 144 ward diffusion process, attacks on DP_{DDPM} and DP_{DDIM}
 145 under the DW_{Fwd}-box setting should yield similar out-
 146 comes. Conversely, from our perspective, which em-
 147 phasizes the importance of stochasticity throughout the
 148 processes, we expect that an attack on DP_{DDIM} with the
 149 DW_{Fwd}-box setting would produce similar results to an
 150 attack on DP_{DDPM} with the DW_{Both}-box setting since
 151 both models would be completely deterministic for the
 152 attacker. To test our hypothesis, we conduct rigorous
 153 evaluations using precise gradients and Expectation over
 154 Transformation (EoT) techniques, as detailed in Sec-
 155 tion 6.1. We employ a 20-step PGD attack with 10 EoT
 156 iterations (l_∞ norm) on the CIFAR-10 dataset. Further
 157 experimental details are discussed in Section 6.2. The
 158 results, depicted in Figure 2, reveal that while DP_{DDPM}
 159 maintains its robustness in the DW_{Fwd}-box setting, DP_{DDIM} significantly loses its robustness. This
 160 finding challenges the traditional view that narrowing the distribution gap during the forward process
 161 enhances robustness. Moreover, in the DW_{Both} setting, the robustness of DP_{DDPM} also deteriorates,
 162 which supports our hypothesis that the robustness of DBP models depends on stochasticity throughout
 163 the DBP processes.

164 Based on these results, we argue that the robustness
 165 of DBP models may primarily stem from the stochas-
 166 ticity throughout the DBP process. Specifically, DBP
 167 models may exploit stochasticity to circumvent the
 168 most effective attack direction. This perspective com-
 169 plements the ideas in previous work on DBP robust-
 170 ness [Xiao et al., 2022; Carlini et al., 2022].

171 4.2 Explaining Stochasticity-Driven Robustness

172 In the previous section, we highlight the sig-
 173 nificance of stochasticity in the robustness of
 174 DBP models. To further investigate the im-
 175 pact of stochasticity on traditional White-box at-
 176 tacks, we analyze the robustness of various DBP
 177 methods such as DiffPure, GDMP, DP_{DDPM}, and
 178 DP_{DDIM} under two conditions: with 10 EoT iter-
 179 ations (EoT10) and without EoT (EoT1). Our re-
 180 sults, detailed in Table 1, show that under White-
 181 box attacks, DBP models retain significant ro-
 182 bustness and that the introduction of EoT leads
 183 to only a marginal reduction. This is in contrast
 184 to the results obtained for DW-box attacks.

185 To further understand robustness performance
 186 under different attack settings, we visualize at-
 187 tack trajectories using t-SNE to map them onto the
 188 xy plane, with corresponding loss values on

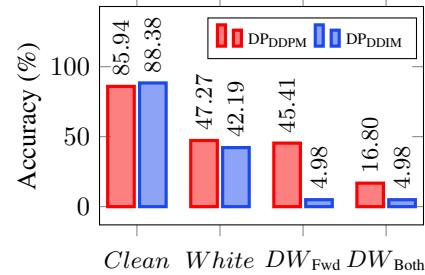


Figure 2: DP_{DDPM} and DP_{DDIM} robustness under different attack settings. Both models lost most of their robustness only when the attacker knows all stochastic elements (DW_{Both}-box for DP_{DDPM} and DW_{Fwd}-box for DP_{DDIM}).

This finding challenges the traditional view that narrowing the distribution gap during the forward process enhances robustness. Moreover, in the DW_{Both} setting, the robustness of DP_{DDPM} also deteriorates, which supports our hypothesis that the robustness of DBP models depends on stochasticity throughout the DBP processes.

Table 1: Evaluation of state-of-the-art DBP methods, EoT significantly influences the evaluation accuracy (%) of model robustness.

	DiffPure	GDMP (MSE)	DP _{DDPM}	DP _{DDIM}
Clean	89.26	91.80	85.94	88.38
PGD20-EoT1	69.04	53.13	60.25	54.59
PGD20-EoT10	55.96	40.97	47.27	42.19

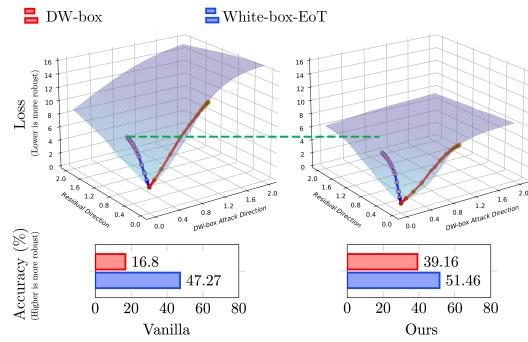


Figure 3: Visualisation of attack trajectories for White-box-EoT attacks and Deterministic White-box attacks on the loss landscape. The loss landscape is not flat in the direction of the Deterministic White-box attack.

189 the z axis. We label the PGD20-EoT10/PGD20-EoT1 attacks within the White-box setting as White-
 190 box-EoT/White-box and the PGD20-EoT1 attack within the DW-box setting as DW-box. As shown
 191 in Figure 1, the trajectories are scattered in all settings, highlighting the stochastic nature of DBP
 192 models. However, DW-box attacks result in a large increase in loss, whereas White-box attacks show
 193 a mild increase. This suggests that DBP models may rely on stochasticity to evade effective attack
 194 directions. In the White-box-EoT setting, trajectories are more concentrated but diverge from the
 195 DW-box attack trajectory, explaining why EoT attacks fail to compromise this defense.

196 From a loss landscape perspective, we suggest that the loss landscape of DBP models is not inherently
 197 smooth. We visualize a White-box-EoT attack and a deterministic White-box attack trajectory on the
 198 loss landscape [Li et al., 2018; Kim et al., 2021]. Experimental details are presented in Appendix D.
 199 As shown in the left side of Figure 3, the trajectory of the White-box-EoT attack diverges from the
 200 Deterministic White-box attack direction, and the loss landscape appears flatter in this direction. This
 201 divergence is possibly due to the failure of the White-box-EoT attack to identify the most effective
 202 direction, supporting the findings from the previous trajectory visualization of Figure 1. Besides, the
 203 significant increase in loss along the direction of the Deterministic White-box attack differs from
 204 the case of AT-trained models, where the landscape may be uniformly flat for all directions [Shafahi
 205 et al., 2019a]. Further evaluation with up to 512 EoT iterations, discussed in Appendix A, confirms
 206 these observations.

207 5 Adversarial Denoising Diffusion Training

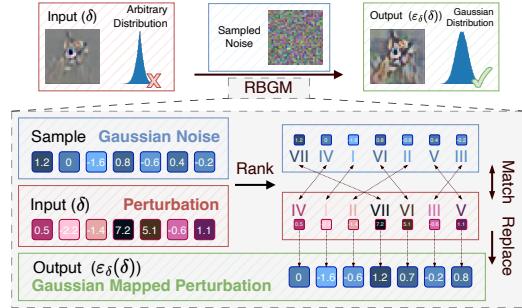
208 In our previous discussions, we note that the robustness of DBP models is primarily due to their
 209 reliance on stochastic elements to evade the most effective attack direction, rather than directly
 210 countering adversarial perturbations like AT models.

211 To improve the robustness of DBP models, we
 212 advocate the integration of adversarial perturba-
 213 tions into the training process. To align these
 214 perturbations with the training of the diffusion
 215 model, we introduce **Rank-Based Gaussian**
 216 **Mapping** (RBGM), a technique designed to
 217 make the perturbations more “Gaussian-like”,
 218 which is elaborated in Section 5.1. Building
 219 on this foundation, we introduce **Adver-
 220 sarial Denoising Diffusion Training** (ADDT).
 221 ADDT uses an iterative two-step approach: the
 222 **Classifier-Guided Perturbation Optimization**
 223 (CGPO) step generates adversarial perturbations,
 224 while the training step trains on these pertur-
 225 bations and updates the diffusion model par-
 226 ameters. This process is illustrated graphically in
 227 Figure 5 and the pseudocode is in Appendix E.
 228 ADDT allows for efficient fine-tuning of clean
 229 diffusion models to proactively counter adver-
 230 sarial perturbations, with a discussion of its efficiency
 presented in Appendix M. Further details on ADDT are discussed in Section 5.

231 5.1 Rank-Based Gaussian Mapping

232 Traditional diffusion models operate under the premise that input images are compromised by
 233 independent Gaussian noise ϵ , as elucidated in Equation (2). This assumption limits our ability
 234 to incorporate perturbations during training that could potentially improve model robustness. To
 235 circumvent this limitation, we propose to employ a modified noise vector, ϵ' , conditioned on the
 236 input x . This modification allows us to tailor ϵ' to possess adversarial properties while maintaining
 237 an approximately Gaussian conditional distribution.

238 To make the perturbations more Gaussian-like, we introduce Rank-Based Gaussian Mapping (RBGM),
 239 illustrated in Figure 4. The RBGM approach, denoted by $\epsilon_\delta()$, takes δ as input and aims to preserve
 240 the relative magnitudes of the elements of δ while ignoring their exact values. The RBGM process
 241 involves sampling a Gaussian tensor, ϵ_s , that matches the dimensions of δ . By replacing the elements



242 Figure 4: Rank-Based Gaussian Mapping. RBGM
 243 trims the input to follow Gaussian distribution. It
 244 samples a Gaussian noise and then replaces ele-
 245 ments in the input with those from the Gaussian
 246 noise, matched according to their respective ranks.

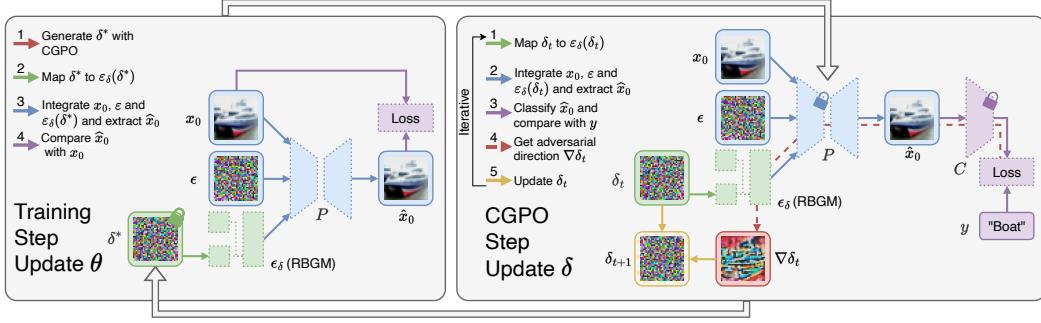


Figure 5: Overview of Adversarial Denoising Diffusion Training (ADDT). ADDT alternates between a CGPO step (right grey box) to refine the perturbations with a frozen diffusion model and classifier, and a training step (left grey box) to update the diffusion model with the refined perturbation. Throughout the process, RBGM is used to make the perturbation more “Gaussian-like”.

242 of δ with the elements of ϵ_s of the same rank, RBGM makes δ more Gaussian-like. To further
 243 augment the Gaussian nature of the noise, we mix the RBGM-mapped perturbation with random
 244 Gaussian noise.

245 By combining the RBGM-induced perturbation with Gaussian noise, we could generate an adversarial
 246 input x'_t as follows:

$$x'_t(x_0, \epsilon, \epsilon_\delta(\delta)) = \sqrt{\alpha_t} x_0 + \sqrt{1 - \lambda_t^2} \sqrt{1 - \bar{\alpha}_t} \epsilon + \lambda_t \sqrt{1 - \bar{\alpha}_t} \epsilon_\delta(\delta), \quad (6)$$

247 where λ_t modulates the level of adversarial perturbations, ensuring that the overall noise remains
 248 largely independent of x_0 and that the perturbations do not overwhelm the denoising model’s learning
 249 capabilities. We determine λ_t using the following formulation:

$$\lambda_t = \text{clip}(\gamma_t \lambda_{\text{unit}}, \lambda_{\min}, \lambda_{\max}), \quad \gamma_t = \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}}, \quad (7)$$

250 where the `clip` function limits λ_t between λ_{\min} and λ_{\max} .

251 5.2 Adversarial Denoising Diffusion Training

252 **Training Step.** The training step is shown to the left of Figure 5. The goal of model optimization is
 253 to subtract the Gaussian noise and the RBGM-mapped perturbations to fully recover the image, this
 254 could be denoted as:

$$\theta^* = \arg \min_{\theta} \mathbb{E} \left[\frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} \|x_0 - P(\theta, x'_t, t)\|_2^2 \right], \quad (8)$$

255 where the expectation is taken with respect to $x_0 \sim \mathcal{D}$, $t \sim \mathcal{U}(\{1, \dots, T\})$, $\epsilon \sim \mathcal{N}(0, I)$. The noise
 256 mixture x'_t is derived from x_0 , ϵ , and δ according to Equation (6), P is a single-step diffusion model,
 257 it takes x'_t and t to predict the original image \hat{x}_0 . Here, instead of adopting multi-step reconstruction
 258 techniques, we use a single-step recovery method. For DDPM and DDIM, this process is governed
 259 by Equation (4). Also note that since the DDPM/DDIM loss is the expectation of the square error of
 260 noise, we add a factor $\sqrt{\bar{\alpha}_t}/\sqrt{1 - \bar{\alpha}_t}$ to keep our loss in line with the traditional DDPM/DDIM loss.

261 **CGPO Step.** The objectives of generative models are distinct from those of downstream perceptual
 262 tasks. To generate perturbations that are meaningful for perceptual tasks, we propose the incorpora-
 263 tion of classifier guidance into the generation of perturbations with Classifier-Guided Perturba-
 264 tion Optimization (CGPO)¹. As shown on the right side of Figure 5, we first reconstruct a clean image
 265 \hat{x}_0 from x'_t with the single-step diffusion model P , we then classify the image with a pre-trained
 266 classifier C . The goal of CGPO is to refine δ to maximize the classification error:

$$\delta^* = \arg \max_{\delta} \mathbb{E} [L(C(P(\theta, x'_t(x_0, \epsilon, \epsilon_\delta(\delta)), t)), y)]. \quad (9)$$

¹Here we use the classifier only for semantic guidance, and the classifier in CGPO doesn’t need to be consistent with the protected model, see Section 6.4.

267 To optimize δ , we employ an iterative process of gradient accumulation, detailed in Appendix E.
 268 Given the non-differentiable nature of RBGM, we accumulate gradients with respect to $\nabla_{\epsilon_\delta(\delta)}$
 269 rather than ∇_δ , which is sufficient for training. In RBGM, we process each channel of every image
 270 independently.

271 6 Experiments

272 6.1 Evaluating Diffusion-Based Purification Robustness

273 Previous assessments of DBP robustness
 274 have often utilized potentially unreliable
 275 methods. In particular, due to the iterative
 276 denoising process in diffusion models,
 277 some studies resort to mathematical ap-
 278 proximations of gradients to reduce mem-
 279 ory constraints [Athalye et al., 2018] or
 280 to circumvent the diffusion process during
 281 backpropagation [Wang et al., 2022]. Fur-
 282 thermore, the reliability of AutoAttack, a
 283 widely used evaluation method, in assess-
 284 ing the robustness of DBP models is ques-
 285 tionable. Although AutoAttack includes a
 286 *Rand* version designed for stochastic mod-
 287 els, Nie et al. have found instances where
 288 the *Rand* version is less effective than the *Standard* version in evaluating DBP robustness [Nie et al.,
 289 2022].

290 To improve the robustness evaluation of diffusion-
 291 based purification (DBP) models, we implement sev-
 292 eral modifications. First, to ensure the accuracy of
 293 the gradient computations, we compute the exact gra-
 294 dient of the entire diffusion classification pipeline. To
 295 mitigate the high memory requirements in the itera-
 296 tive denoising steps, we use gradient checkpointing
 297 techniques to optimize memory usage. In addition,
 298 to deal with the stochastic nature of the DBP process,
 299 we incorporate the Expectation over Transformation
 300 (EoT) method to average gradients across different
 301 attacks. We adopt EoT with 10 iterations, and a de-
 302 tails discussion of the choice of EoT iterations can
 303 be found in Appendix A. We also use the Projected
 304 Gradient Descent (PGD) attack instead of AutoAt-
 305 tack for our evaluations². Our revised robustness
 306 evaluation revealed that DBP models, such as DiffPure and GDMP, perform worse than originally
 307 claimed. DiffPure’s accuracy dropped from a claimed 70.64% to an actual 55.96%, and GDMP’s
 308 from 90.10% to 40.97%. These results emphasize the urgent need for more accurate and reliable
 309 evaluation methods to properly assess the robustness of DBP models.

310 6.2 Experiment Setup

311 **Classifier.** We train a WideResNet-28-10 classifier for 200 epochs following the methods in [Yoon
 312 et al., 2021; Wang et al., 2022] and achieve 95.12% accuracy on CIFAR-10 and 76.66% on CIFAR-
 313 100 dataset.

314 **DBP.** For DBP forward process timestep, DiffPure employs a continuous-time VPSDE (DDPM++)
 315 model, selecting $t^* = 0.1$. For discrete timesteps in Equation (2) of DP_{DDPM} and DP_{DDIM}, we set
 316 $t^* = 0.1 \times T$. We also implement DP_{EDM}. The details are discussed in Appendix F.

Table 2: Clean and robust accuracy (%) on CIFAR-10 obtained by different DBP methods. All methods show consistent improvement fine-tuned with ADDT.

Diffusion model	DBP model	<i>Clean</i>	l_∞	l_2
-	-	95.12	0.00	1.46
DDIM	DP _{DDIM}	88.38	42.19	70.02
	DP _{DDIM+Ours}	88.77	46.48	71.19
	GDMP (No Guided) [Wang et al., 2022]	91.41	40.82	69.63
	GDMP (MSE) [Wang et al., 2022]	91.80	40.97	70.02
DDPM	GDMP (SSIM) [Wang et al., 2022]	92.19	38.18	68.95
	DP _{DDPM}	85.94	47.27	69.34
	DP _{DDPM+Ours}	85.64	51.46	70.12
	COUP [Zhang et al., 2024]	90.33	50.78	71.19
DDPM++	DiffPure	89.26	55.96	75.78
	DiffPure+Ours	89.94	62.11	76.66
	EDM	86.43	62.50	76.86
	DP _{EDM} (Appendix F)	86.33	66.41	79.16
	DP _{EDM+Ours} (Appendix F)			

Table 3: Clean and robust accuracy (%) on DP_{DDPM}. ADDT improve robustness across different NFE, especially at lower NFE (*: default DDPM generation setting; -: classifier only).

Dataset	NFE	Vanilla			Ours		
		<i>Clean</i>	l_∞	l_2	<i>Clean</i>	l_∞	l_2
CIFAR-10	-	95.12	0.00	1.46	95.12	0.00	1.46
	5	49.51	21.78	36.13	59.96	30.27	41.99
	10	73.34	36.72	55.47	78.91	43.07	62.97
	20	81.45	45.21	65.23	83.89	48.44	69.82
	50	85.54	46.78	68.85	85.45	50.20	69.04
	100*	85.94	47.27	69.34	85.64	51.46	70.12
CIFAR-100	-	76.66	0.00	2.44	76.66	0.00	2.44
	5	17.29	3.71	9.28	21.78	6.25	13.77
	10	34.08	10.55	19.24	40.62	14.55	27.25
	20	48.05	17.68	30.66	53.32	18.65	36.13
	50	55.57	20.02	37.70	59.47	22.75	40.72
	100*	57.52	20.41	37.89	59.18	23.73	41.70

²We discover a bug in the *Rand* version of AutoAttack that causes it to overestimate the robustness of DBP. After fixing this, AutoAttack gives similar results to PGD attacks, but at a much higher computational cost. We discuss this in detail in Appendix J.

317 **ADDT.** We fine-tune our diffusion model based on the CIFAR-
 318 10 pre-trained exponential moving
 319 average (EMA) clean model
 320 closely following its training set-
 321 ting [Ho et al., 2020] (trans-
 322 formed into Huggingface Dif-
 323 fusers format by [Fang et al.,
 324 2023]). The CIFAR-100 clean
 325 model is fine-tuned from the
 326 CIFAR-10 clean model for 100
 327 epochs. ADDT models are
 328 fine-tuned from clean pre-trained
 329 models for 100 epochs, with
 330 guidance from the WideResNet-28-10 classifier. In CGPO we adopt $\lambda_{unit} = 0.03$, $\lambda_{min} = 0$, $\lambda_{max} = 0.3$ and refine δ for 5 iterations. Fine-tuning from clean models with ADDT is quite
 331 efficient, taking 12 hours on four NVIDIA GeForce RTX 2080ti GPUs for DP_{DDPM}/DP_{DDIM} models.
 332

333 **Robustness Evaluation.** We adopt PGD20-EoT10 attack for robustness evaluation. For l_∞ attacks
 334 we set $\alpha = 2/255$ and $\epsilon = 8/255$; for l_2 attacks we set $\alpha = 0.1$ and $\epsilon = 0.5$. Due to computational
 335 constraints, we test the first 1024 images from the CIFAR-10/100 datasets. Note that it still takes 5
 336 hours on four NVIDIA GeForce RTX 2080ti GPUs to test DP_{DDPM}/DP_{DDIM}.
 337

338 6.3 Comparison with State-of-the-Art Approaches

339 We apply ADDT fine-tuning to a set of diffusion models and apply DiffPure-style DBP with the
 340 refined models. We then compare their performance with state-of-the-art DBP methods on the
 341 CIFAR-10 dataset. The outcomes, detailed in Table 2, reveal that ADDT fine-tuning enhances the
 342 robustness of these models, enabling them to reach state-of-the-art performance.

343 6.4 Defense Performance in Extensive Scenarios

344 **Performance on Different Classifiers.** We evaluate the cross model protection ability of ADDT
 345 fine-tuned models, as shown in Table 4. The results indicate that fine-tuning with WRN-28-10 guided
 346 ADDT could enhance the ability to
 347 protect different classifiers. Notably,
 348 using DP_{EDM}, we achieve 69.63%
 349 l_∞ robust accuracy on a WRN-70-
 350 16 classifier. This demonstrates that
 351 our method, without classifier-specific
 352 fine-tuning, can achieve comparable
 353 results to state-of-the-art AT-based
 354 models.

355 **Performance under Acceleration.**
 356 Speeding up the diffusion process by
 357 omitting intermediate steps has be-
 358 come common practice in the use of
 359 diffusion models [Nichol and Dhari-
 360 wal, 2021; Song et al., 2020a]. Here
 361 we evaluate the robustness of accelerated DBP models. Acceleration is measured by the number
 362 of neural function evaluations (NFE), which indicates the number of evaluation steps performed
 363 during the DBP backtracking process. For our experiments, we set $t^* = 0.1 \times T$ and accelerate the
 364 process by excluding intermediate time steps. For example, with an NFE of 5, the time steps for
 365 the DBP backward process would be $t = [100, 80, 60, 40, 20, 0]$. The results are detailed in Table 3.
 366 Our method improves the robustness of both DP_{DDPM} models. Note that the performance of DP_{DDPM}
 367 varies significantly between different NFE. This may be due to the fact that DDPM introduces stochas-
 368 ticity (Gaussian noise) at each reverse step; with fewer reverse steps, it's stochasticity decreases.

Table 4: Clean and robust accuracy (%) on CIFAR-10, obtained by different classifiers. ADDT (WRN-28-10 guidance) improves robustness in protecting different subsequent classifiers. (*: the classifier used in ADDT fine-tuning).

Method	Classifier	Vanilla			Ours		
		Clean	l_∞	l_2	Clean	l_∞	l_2
DP _{DDPM} -1000	VGG-16 [Simonyan and Zisserman, 2014]	84.77	41.99	66.89	85.06	46.09	67.87
	ResNet-50 [He et al., 2016]	83.11	44.04	67.58	83.84	48.14	67.87
	WRN-28-10* [Zagoruyko and Komodakis, 2016]	85.94	47.27	69.34	85.64	51.46	70.12
	WRN-70-16 [Zagoruyko and Komodakis, 2016]	88.43	48.93	70.31	87.84	52.54	70.70
DP _{DDIM} -100	ViT-B [Dosovitskiy et al., 2020]	85.45	45.61	69.53	85.25	48.63	69.92
	VGG-16 [Simonyan and Zisserman, 2014]	87.16	29.00	61.82	87.55	35.06	66.11
	ResNet-50 [He et al., 2016]	86.04	31.74	62.11	86.57	38.77	65.82
	WRN-28-10* [Zagoruyko and Komodakis, 2016]	88.96	43.16	67.58	88.18	47.85	70.61
DP _{EDM}	WRN-70-16 [Zagoruyko and Komodakis, 2016]	84.40	39.16	68.36	84.96	47.66	69.14
	ViT-B [Dosovitskiy et al., 2020]	88.77	34.38	65.72	88.48	41.02	68.65
	WRN-28-10* [Zagoruyko and Komodakis, 2016]	86.43	62.50	76.86	86.33	66.41	79.16
	WRN-70-16 [Zagoruyko and Komodakis, 2016]	86.62	65.62	76.46	86.43	69.63	78.91

Table 5: Clean and robust accuracy (%) on CIFAR-10 fine-tuned with different methods. Vanilla fine-tuning and fine-tuning with non-classifier-guided perturbations did not improve robustness.

Method	NFE	Vanilla Fine-tune Only			MSE Distance Guided		
		Clean	l_∞	l_2	Clean	l_∞	l_2
DP _{DDPM}	5	47.27	21.88	33.11	49.32	21.19	35.74
	10	71.58	34.77	52.64	73.34	36.43	55.96
	20	81.93	42.29	63.77	83.79	42.87	66.89
	50	84.18	47.56	65.53	85.16	47.27	67.87
	100*	85.25	47.27	68.26	86.91	46.97	70.80
DP _{DDIM}	5	89.26	41.11	67.87	89.45	39.36	67.58
	10*	88.87	41.41	67.19	89.36	40.92	67.68
	20	89.45	40.23	67.77	88.48	42.68	68.55
	50	88.77	42.09	67.38	88.18	41.02	68.16
	100	88.18	40.62	68.26	89.26	40.33	68.36

Table 6: Clean and robust accuracy (%) on Tiny-ImageNet with WRN-28-10 classifier. ADDT improve DBP robustness on Tiny-ImageNet (-: classifier only).

Method	Vanilla			Ours		
	Clean	l_∞	l_2	Clean	l_∞	l_2
-	71.37	0.00	0.00	-	-	-
DP _{DDPM} -1000	57.13	11.82	46.68	56.15	13.57	48.54
DP _{DDIM} -100	60.35	4.79	39.75	60.45	5.86	40.82
DP _{PEDM}	57.03	19.14	46.00	56.45	20.61	47.95

369 Additionally, the generation capability of DDPM is sensitive to the skipping of intermediate steps.
370 We also conducted an evaluation of DP_{DDIM} models, as detailed in Appendix G.

371 **Performance on Tiny-ImageNet.** We test DBP robustness on Tiny-ImageNet dataset [Le and Yang,
372 2015] in Table 6, with training and evaluation settings following CIFAR-10. The result shows that
373 ADDT improves robustness on Tiny-ImageNet dataset.

374 6.5 Ablation Study and Analysis

375 **RBGM.** We compare the generative ability of diffusion models fine-tuned with RBGM-mapped and
376 l_∞ perturbations by comparing their Fréchet Inception Distance (FID) scores, as shown in Table 7. The
377 results show that diffusion models fine-tuned with RBGM-mapped perturbations maintain generation
378 quality comparable to the vanilla diffusion model, while models fine-tuned with l_∞ perturbations
379 show degraded performance. We also observe that training with RBGM-mapped perturbations
380 generalized better to different attacks. Experimental details and additional tests are presented in
381 Appendix I.

382 **CGPO.** We examine the effect of vanilla fine-
383 tuning and ADDT fine-tuning with MSE dis-
384 tance guidance (without classifier guidance) in
385 Table 5. Neither method improves the robust-
386 ness of DBP models.

387 **Revisiting DBP Robustness.** We re-examine
388 robustness under the Deterministic White-box
389 setting by evaluating the performance of diffu-
390 sion models after ADDT fine-tuning, as shown
391 in Figure 6. The fine-tuned models show signif-
392 icantly higher robustness under the Deterministic
393 White-box setting, indicating that our method
394 makes the attack more difficult without relying
395 on the evasion effect of stochastic elements. We
396 provide a more detailed experiment considering
397 different NFE in Appendix K. We also compare
398 the loss landscapes of ADDT fine-tuned mod-
399 els and vanilla diffusion models, as shown in Figure 3. This comparison shows that our method
400 effectively smooths the loss landscape of DBP models.

401 7 Conclusion

402 This study offers a new perspective on the robustness of Diffusion-Based Purification (DBP) models,
403 emphasizing the crucial role of stochasticity and challenging the traditional view that robustness
404 is mainly derived from minimizing the distribution gap through the forward diffusion process. We
405 introduce a Deterministic white-box (DW-box) attack scenario and show that DBP models are based
406 on stochastic elements to evade effective attack directions. To further enhance the robustness of DBP
407 models, we have developed Adversarial Denoising Diffusion Training (ADDT) and Rank-Based
408 Gaussian Mapping (RBGM). ADDT integrates adversarial perturbations into the training process,
409 while RBGM trims perturbations to more closely resemble Gaussian distributions. Our empirical
410 results confirm that ADDT achieves robustness improvements of up to 6% over conventional DBP
411 models, underscoring the effectiveness of our proposed enhancements.

Table 7: FID score of DDPM for CIFAR-10 fine-tuned to different perturbations (the lower the better). Fine-tuning with RBGM-mapped perturbations yields lower FID scores than l_∞ perturbations.

	Vanilla	Clean	Fine-tune	Ours	Ours _{l_∞}
FID	3.196	3.500	5.190	13.608	

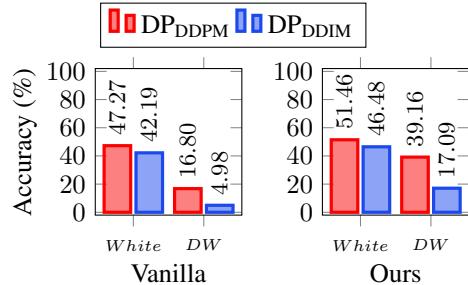


Figure 6: Revisiting robustness under Deterministic White-box setting. ADDT improves robustness under both White-box and Deterministic White-box setting, implying that ADDT strengthens the models' ability to counter adversarial inputs.

412 **References**

- 413 Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of
414 security: Circumventing defenses to adversarial examples. In *International conference on machine*
415 *learning*, pages 274–283. PMLR, 2018.
- 416 Nicholas Carlini, Florian Tramer, Krishnamurthy Dj Dvijotham, Leslie Rice, Mingjie Sun, and J Zico
417 Kolter. (certified!!) adversarial robustness for free! In *The Eleventh International Conference on*
418 *Learning Representations*, 2022.
- 419 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of
420 diverse parameter-free attacks. In *International conference on machine learning*, pages 2206–2216.
421 PMLR, 2020.
- 422 Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo Debenedetti, Nicolas Flam-
423 marion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial
424 robustness benchmark. *arXiv preprint arXiv:2010.09670*, 2020.
- 425 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
426 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
427 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
428 *arXiv:2010.11929*, 2020.
- 429 Gongfan Fang, Xinyin Ma, and Xinchao Wang. Structural pruning for diffusion models. *arXiv*
430 *preprint arXiv:2305.10924*, 2023.
- 431 Yue Gao, Ilia Shumailov, Kassem Fawaz, and Nicolas Papernot. On the limitations of stochastic
432 pre-processing defenses. *Advances in Neural Information Processing Systems*, 35:24280–24294,
433 2022.
- 434 Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
435 examples. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego,*
436 *CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- 437 Sven Gowal, Chongli Qin, Jonathan Uesato, Timothy A. Mann, and Pushmeet Kohli. Uncovering the
438 limits of adversarial training against norm-bounded adversarial examples. *CoRR*, abs/2010.03593,
439 2020.
- 440 Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and
441 Timothy A Mann. Improving robustness using generated data. *Advances in Neural Information*
442 *Processing Systems*, 34:4218–4233, 2021.
- 443 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
444 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
445 pages 770–778, 2016.
- 446 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint*
447 *arxiv:2006.11239*, 2020.
- 448 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
449 based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577,
450 2022.
- 451 Hoki Kim, Woojin Lee, and Jaewook Lee. Understanding catastrophic overfitting in single-step
452 adversarial training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages
453 8119–8127, 2021.
- 454 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 455 Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape
456 of neural nets. *Advances in neural information processing systems*, 31, 2018.

- 457 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
458 Towards deep learning models resistant to adversarial attacks. In *6th International Conference on*
459 *Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference*
460 *Track Proceedings*. OpenReview.net, 2018.
- 461 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
462 In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- 463 Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar.
464 Diffusion models for adversarial purification. In *International Conference on Machine Learning*
465 (*ICML*), 2022.
- 466 OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.
- 467 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.
468 Robust speech recognition via large-scale weak supervision. *ArXiv*, abs/2212.04356, 2022.
- 469 Sylvestre-Alvise Rebiffé, Sven Gowal, Dan Andrei Calian, Florian Stimberg, Olivia Wiles, and Timo-
470 thy A. Mann. Fixing data augmentation to improve adversarial robustness. *ArXiv*, abs/2103.01946,
471 2021.
- 472 Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-gan: Protecting classifiers against
473 adversarial attacks using generative models. *arXiv preprint arXiv:1805.06605*, 2018.
- 474 Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John P. Dickerson, Christoph Studer, Larry S.
475 Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *Neural Information*
476 *Processing Systems*, 2019a.
- 477 Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer,
478 Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! *Advances in*
479 *Neural Information Processing Systems*, 32, 2019b.
- 480 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
481 recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 482 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models.
483 *arXiv:2010.02502*, 2020a.
- 484 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
485 *Advances in neural information processing systems*, 32, 2019.
- 486 Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Lever-
487 aging generative models to understand and defend against adversarial examples. In *International*
488 *Conference on Learning Representations*, 2018.
- 489 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
490 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
491 *arXiv:2011.13456*, 2020b.
- 492 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
493 Poole. Score-based generative modeling through stochastic differential equations. In *International*
494 *Conference on Learning Representations*, 2021.
- 495 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow,
496 and Rob Fergus. Intriguing properties of neural networks. In *2nd International Conference on*
497 *Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track*
498 *Proceedings*, 2014.
- 499 Jinyi Wang, Zhaoyang Lyu, Dahua Lin, Bo Dai, and Hongfei Fu. Guided diffusion model for
500 adversarial purification. *arXiv preprint arXiv:2205.14969*, 2022.
- 501 Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion
502 models further improve adversarial training. In *International Conference on Machine Learning*
503 (*ICML*), 2023.

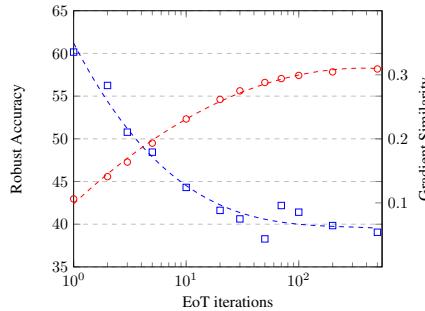
- 504 Quanlin Wu, Hang Ye, and Yuntian Gu. Guided diffusion model for adversarial purification from
505 random noise. *arXiv preprint arXiv:2206.10875*, 2022.
- 506 Chaowei Xiao, Zhongzhu Chen, Kun Jin, Jiongxiao Wang, Weili Nie, Mingyan Liu, Anima Anand-
507 kumar, Bo Li, and Dawn Song. Densepure: Understanding diffusion models towards adversarial
508 robustness. *arXiv preprint arXiv:2211.00322*, 2022.
- 509 Jongmin Yoon, Sung Ju Hwang, and Juho Lee. Adversarial purification with score-based generative
510 models. In *International Conference on Machine Learning*, pages 12062–12072. PMLR, 2021.
- 511 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*,
512 2016.
- 513 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
514 Theoretically principled trade-off between robustness and accuracy. In *International conference on
515 machine learning*, pages 7472–7482. PMLR, 2019.
- 516 Mingkun Zhang, Jianing Li, Wei Chen, Jiafeng Guo, and Xueqi Cheng. Classifier guidance enhances
517 diffusion-based adversarial purification by preserving predictive information, 2024.

518 **A Influence of EoT Iterations on DBP Robustness Evaluation**

519 In this section, we examine how the number of EoT iterations influences the DBP robustness
 520 evaluation. As previously discussed in Section 4.1, the Deterministic White-box attack could find
 521 the most effective attack direction. To quantify the impact of EoT iterations, we compare the attack
 522 direction of the standard White-box-EoT across various numbers of EoT iterations with that of the
 523 Deterministic White-box.

524 See Figure 7 for a visual explanation, where the red line shows the DBP accuracy after attack, and the
 525 blue line shows the similarity between the attack directions of the White-box-EoT and Deterministic
 526 White-box. The trend is clear: more EoT iterations lead to greater similarity and lower model
 527 accuracy, the rate of increase in similarity and the rate of decrease in accuracy both tend to slow down
 528 with further iterations.

529 Balancing computational cost and evaluation accuracy, we chose the PGD20-EoT10 configuration
 for our robustness evaluation.



530 Figure 7: Robust accuracy and gradient similarity on DP_{DDPM} for CIFAR-10, obtained by different
 531 EoT iterations. As the number of EoT iterations increases, the gradient similarity between the
 532 White-box-EoT attack direction and the Deterministic White-box attack direction increases and the
 533 robust accuracy decreases.

534

535 **B Available Information in Different Attack Settings**

536 We provide information available to attackers in different attack settings in Table 8.

537 Table 8: Available information in different attack settings. ϵ is the randomly generated Gaussian noise
 538 added in DBP diffusion forward and backward process, DDIM does not have ϵ in DBP diffusion
 539 backward process.

Attacker Known Things\Attack Settings	White-box	DW _{Fwd} -box	DW _{Both} -box
Model Architecture and Parameters	✓	✓	✓
Input Images, Class Label	✓	✓	✓
DDPM/DDIM ϵ in Equation (2)	✗	✓	✓
DDPM ϵ in Equation (5)	✗	✗	✓

538

539 **C DBP Models Employing Different Stochastic Elements Cannot Be Attacked
 540 All at Once**

541 Previous research has questioned whether stochasticity can improve robustness, arguing that it can
 542 produce obfuscated gradients that give a false sense of security [Athalye et al., 2018]. To investigate
 543 this, we implement DW_{Semi}-box, a semi-stochastic setting that restricts the stochastic elements to a
 544 limited set of options. Our results show that stochasticity can indeed improve robustness, even when
 545 the attacker has full knowledge of all the possible options for stochastic elements. For a detailed
 546 analysis of our results, see Appendix C. The potential to further improve robustness by increasing the
 547 stochasticity of the model is discussed in Appendix H.

542 Building on the concept of Deterministic White-box, we further propose $DW_{\text{semi-128}}$ to explore whether
 543 stochasticity can indeed improve robustness. Unlike under Deterministic White-box, where the
 544 attacker attacks a DBP model under the exact set of stochastic noise used in the evaluation, $DW_{\text{semi-128}}$
 545 relaxes the stochastic elements to a limited set of options, the attacker should simultaneously attack
 546 over 128 different sets of stochastic noise. It uses the average adversarial direction from these 128
 547 noise settings (EoT-128) to perturb the DBP model. To understand the impact of stochasticity, we
 548 analyze the changes of the model loss under DW -box attack and $DW_{\text{semi-128}}$ attack. We plot these
 549 changes by adjusting a factor k to modify an image x with a perturbation σ , evaluating the loss
 550 at $x + k\sigma$ where k varies from -16 to 16 . We generate perturbations with l_∞ Fast Gradient Sign
 551 Method (FGSM) [Goodfellow et al., 2015] with magnitude $1/255$. The plot is evaluated using
 552 WideResNet-28-10 with DP_{DDPM} over the first 128 images of CIFAR-10 dataset.

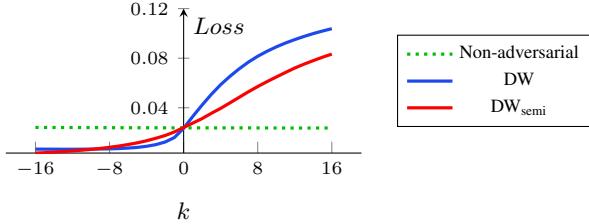


Figure 8: Impact of stochasticity on perturbation efficacy. Perturbations created under DW_{semi} -box setting are less potent compared to DW -box setting. For non-adversarial perturbations, we randomly assign each element a value of either $1/255$ or $-1/255$.

553 As Figure 8 shows, in the Deterministic White-box setting, the perturbations significantly increase
 554 the loss, proving their effectiveness. However, for $DW_{\text{semi-128}}$, where the attack spans multiple noise
 555 setting, the increase in loss is more moderate. This suggests that even when the attackers are fully
 556 informed about the stochastic noise choices, stochasticity still improves the robustness of the DBP.
 557 This challenges the notion that there exists a vulnerable direction that is effective for all stochastic
 558 noise.

559 D Experimental Setting of Visualization of the Attack Trajectory

560 We visualize the attack by plotting the loss landscape and trace the trajectories of EoT attack under
 561 White-box setting and the Deterministic White-box setting in Figure 3. We run a vanilla PGD20-
 562 EoT10 attack under White-box setting and a PGD20 attack under Deterministic White-box setting.
 563 We then expand a 2D space using the final perturbations from these two attacks, draw the loss
 564 landscape, and plot the attack trajectories on it. Note that the two adversarial perturbation directions
 565 are not strictly orthogonal. To extend this 2D space, we use the Deterministic White-box attack
 566 direction and the orthogonal component of the EoT attack direction. Note that the endpoints of both
 567 trajectories lie exactly on the loss landscape, while intermediate points are projected onto it. The plot
 568 is evaluated using WideResNet-28-10 with DP_{DDPM} over the first 128 images of CIFAR-10 dataset.

569 E Pseudo-code of ADDT

570 The pseudo-code for adopting ADDT within DDPM and DDIM framework is shown in Algorithm 1.

571 F Adopting VPSDE(DDPM++) and EDM Models in DBP

572 In the previous discussion of the robustness of DBP models, as detailed in Section 4.1, our focus
 573 was primarily on the DDPM and DDIM models. We now extend our analysis to include VPSDE
 574 (DDPM++) and EDM [Karras et al., 2022] models. VPSDE (DDPM++) is the diffusion model used
 575 in DiffPure.

576 From a unified perspective, diffusion processes can be modeled by stochastic differential equations
 577 (SDE) [Song et al., 2021]. The forward SDE, as described in Equation (10), converts a complex
 578 initial data distribution into a simpler, predetermined prior distribution by progressively infusing

Algorithm 1 Adversarial Denoising Diffusion Training (ADDT)

Require: \mathbf{x}_0 is image from training dataset, y is the class label of the image, C is the classifier, \mathbf{P} is one-step diffusion reverse process and θ is it's parameter, L is CrossEntropy Loss.

- 1: **for** \mathbf{x}_0, y in the training dataset **do**
- 2: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 3: $\lambda_t = \text{clip}(\gamma_t \lambda_{\text{unit}}, \lambda_{\min}, \lambda_{\max})$, where $\gamma_t = \frac{\sqrt{\alpha_t}}{\sqrt{1-\alpha_t}}$
- 4: Init δ to a small random vector.
- 5: **for** 1 to ADDT_{iterations} **do**
- 6: $\epsilon \sim \mathcal{N}(0, I)$
- 7: $\epsilon' = \text{RBGM}(\delta, \epsilon)$
- 8: $\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \lambda_t^2} \sqrt{1 - \alpha_t} \epsilon + \lambda_t \sqrt{1 - \alpha_t} \epsilon'$
- 9: $\delta = \delta + \nabla_{\epsilon'} L(C(\mathbf{P}(\mathbf{x}_t, t), y))$
- 10: **end for**
- 11: $\epsilon \sim \mathcal{N}(0, I)$
- 12: $\epsilon' = \text{RBGM}(\delta, \epsilon)$
- 13: $\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \lambda_t^2} \sqrt{1 - \alpha_t} \epsilon + \lambda_t \sqrt{1 - \alpha_t} \epsilon'$
- 14: Take a gradient descent step on:

$$\nabla_{\theta} \left\| \frac{\sqrt{\alpha_t}}{\sqrt{1-\alpha_t}} (\mathbf{x}_0 - \mathbf{P}(\mathbf{x}_t, t)) \right\|_2^2$$
- 15: **end for**

Diffusion Unet ϵ_{θ} predicts the Gaussian noise added to the image, adopting Equation (4) in the paper, we have $\mathbf{P}(\mathbf{x}_t, t) = (\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_{\theta}(\mathbf{x}_t, t)) / \sqrt{\alpha_t}$

579 noise. This can also be done in a single step, as shown in Equation (11), mirroring the strategy of
580 DDPM described in Equation (2). Reverse SDE, as explained in Equation (12), reverses this process,
581 restoring the noise distribution to the original data distribution, thus completing the diffusion cycle.

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}, \quad (10)$$

$$p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) = \mathcal{N}\left(\mathbf{x}(t); e^{-\frac{1}{4}t^2(\bar{\beta}_{\max}-\bar{\beta}_{\min})-\frac{1}{2}t\bar{\beta}_{\min}} \mathbf{x}(0), \mathbf{I} - \mathbf{I} e^{-\frac{1}{2}t^2(\bar{\beta}_{\max}-\bar{\beta}_{\min})-t\bar{\beta}_{\min}}\right), \quad t \in [0, 1] \quad (11)$$

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g^2(t)\nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t)d\bar{\mathbf{w}}. \quad (12)$$

584 The reverse process of SDEs also derives equivalent ODEs Equation (13) for fast sampling and exact
585 likelihood computation, and this Score ODEs corresponds to DDIM.

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g(t)^2\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt. \quad (13)$$

586 By modulating the stochasticity, we can craft a spectrum of semi-stochastic models that bridge pure
587 SDEs and deterministic ODEs, offering a range of stochastic behaviors.

588 EDM provides a unified framework to synthesize the design principles of different diffusion models
589 (DDPM, DDIM, iDDPM [Nichol and Dhariwal, 2021], VPSDE, VESDE [Song et al., 2021]). Within
590 this framework, EDM incorporates efficient sampling methods, such as the Heun sampler, and
591 introduces optimized scheduling functions $\sigma(t)$ and $s(t)$. This allows EDM to achieve state-of-the-art
592 performance in generative tasks.

593 EDM forward process could be presented as:

$$\mathbf{x}_t = \mathbf{x}_0 + \sigma(t^*) * \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}), \quad (14)$$

594 where we choose $\sigma(t^*) = 0.5$ for clean and robust accuracy tradeoff. And for reverse process, EDM
595 incorporates a parameter S_{churn} to modulate the stochastic noise infused during the reverse process.
596 For our experiments, we choose 50 reverse steps (50 NFE, NFE is Function of Neural Function
597 Evaluations), configured the parameters with $S_{min} = 0.01$, $S_{max} = 0.46$, $S_{noise} = 1.007$, and
598 designate $S_{churn} = 0$ to represent EDM-ODE, $S_{churn} = 6$ to represent EDM-SDE.

599 As shown in Table 9, our ADDT could also increase the robustness of DP_{EDM}.

600 G ADDT Results on DP_{DDIM}

601 As shown in Table 10, the performance of DP_{DDIM} is less sensitive to the number of function
602 evaluations (NFE). Additionally, ADDT consistently improved the robustness of DP_{DDIM}.

Table 9: Clean and robust accuracy on DP_{EDM} for CIFAR-10. ADDT improves robustness in both DP_{EDM-SDE} and DP_{EDM-ODE}.

Type	Vanilla		Ours	
	DP _{EDM-SDE}	DP _{EDM-ODE}	DP _{EDM-SDE}	DP _{EDM-ODE}
Clean	86.43	87.99	86.33	87.99
l_∞	62.50	60.45	66.41	64.16
l_2	76.86	75.49	79.16	77.15

Table 10: Clean and robust accuracy (%) on DP_{DDIM}. ADDT improve robustness across different NFE (*: default DDIM generation setting, -: classifier only).

Dataset	NFE	Vanilla			Ours		
		Clean	l_∞	l_2	Clean	l_∞	l_2
CIFAR-10	-	95.12	0.00	1.46	95.12	0.00	1.46
	5	89.65	42.19	68.65	88.57	47.27	70.61
	10*	88.96	43.16	67.58	88.18	47.85	70.61
	20	87.89	41.70	69.24	88.67	48.63	69.73
	50	88.96	42.48	68.85	88.57	46.68	69.24
	100	88.38	42.19	70.02	88.77	46.48	71.19
CIFAR-100	-	76.66	0.00	2.44	76.66	0.00	2.44
	5	62.11	15.43	35.74	62.79	17.58	38.87
	10*	62.21	15.33	36.52	64.45	20.02	39.26
	20	63.67	15.62	37.89	65.23	18.65	40.62
	50	62.40	16.31	37.79	63.87	19.14	39.94
	100	63.28	15.23	36.62	66.02	18.85	39.84

603 H Strengthening DBP via Augmented Stochasticity

604 Song *et al.* present a Predictor-Corrector sampler for SDEs reverse process for VPSDE (DDPM++)
605 (as detailed in Appendix G of [Song et al., 2021]). However, standard implementations of VPSDE
606 (DDPM++) typically use only the Predictor. Given our hypothesis that stochasticity contributes to
607 robustness, we expect that integrating the Corrector sampler into VPSDE (DDPM++) would further
608 enhance the robustness of DBP models. Our empirical results, as shown in Table 11, confirm that the
609 inclusion of a Corrector to VPSDE (DDPM++) indeed improve the model’s defenses ability against
610 adversarial attacks with l_∞ norm constraints. This finding supports our claim that the increased
611 stochasticity can further strengthen DBP robustness. Adding Corrector is also consistent with ADDT.
612 Note that the robustness against l_2 norm attacks does not show a significant improvement with the
613 integration of the Extra Corrector. A plausible explanation for this could be that the robustness under
614 l_2 attacks is already quite strong, and the compromised performance on clean data counteracts the
increase in robustness.

Table 11: Clean and robust accuracy on DP_{DDPM++} for CIFAR-10. Both extra Corrector and ADDT fine-tuning improved robustness.

Type	Vanilla	Extra Corrector	ADDT	ADDT+Extra Corrector
Clean	89.26	85.25	89.94	85.55
l_∞	55.96	59.77	62.11	65.23
l_2	75.78	74.22	76.66	76.66

615

616 I Evaluating RBGM-Mapped Perturbations

617 In Section 6.5, we briefly explore the generation capabilities of diffusion models trained with
618 RBGM-mapped and l_∞ perturbations. Here, we provide details of the experiment and delve deeper
619 into their robustness. To train with l_∞ perturbations, we adjust ADDT, replacing RBGM-mapped
620 perturbations with l_∞ perturbations. Here, instead of converting accumulated gradients to Gaussian-
621 like perturbations, we use a 5-step projected gradient descent (PGD-5) approach. We also set
622 $\lambda_{unit} = 1$, $\lambda_{min} = 0$, $\lambda_{max} = 10$. We refer to this modified training protocol as ADDT _{l_∞} .

Table 12: Clean and robust accuracy on DBP models trained with different perturbations for CIFAR-10. While ADDT simultaneously improves clean accuracy and robustness against both l_2 and l_∞ attacks. ADDT $_{l_\infty}$ primarily improves performance against l_∞ attacks.

Method	Dataset	Vanilla			ADDT			ADDT $_{l_\infty}$		
		Clean	l_∞	l_2	Clean	l_∞	l_2	Clean	l_∞	l_2
DPDDPM-1000	CIFAR-10	85.94	47.27	69.34	85.64	51.46	70.12	84.47	52.64	68.55
	CIFAR-100	57.52	20.41	37.89	59.18	23.73	41.70	57.81	23.24	40.04
DPDDIM-100	CIFAR-10	88.38	42.19	70.02	88.77	46.48	71.19	88.48	50.49	70.31
	CIFAR-100	63.28	15.23	36.62	66.02	18.85	39.84	64.84	20.31	39.36

We evaluate the clean and robust accuracy of ADDT and ADDT $_{l_\infty}$ fine-tuned models. These models exhibit different behaviors. As shown in Table 12, while Gaussian-mapped perturbations can simultaneously improve clean accuracy and robustness against both l_2 and l_∞ attacks, training with l_∞ perturbations primarily improves performance against l_∞ attacks.

J Evaluating under Fixed AutoAttack

AutoAttack [Croce and Hein, 2020], an ensemble of White-box and Black-box attacks, is a popular benchmark for evaluating model robustness. It is used in RobustBench [Croce et al., 2020] to evaluate over 120 models. However, Nie et al. [Nie et al., 2022] found that the *Rand* version of AutoAttack, designed to evaluate stochastic defenses, sometimes yields higher accuracy than the *Standard* version, intended for deterministic methods. Our comparison of AutoAttack and PGD20-EoT10 in Table 13 also shows that the *Rand* version of AutoAttack gives higher accuracy than the PGD20-EoT10 attack.

We attribute this to the sample selection of AutoAttack. As an ensemble of attack methods, AutoAttack selects the final adversarial sample from either the original input or the attack results. However, the original implementation neglects stochasticity and considers a adversarial sample to be sufficiently adversarial if it gives a false result in one evaluation. To fix this, we propose a 20-iteration evaluation and selects the adversarial example with the lowest accuracy. The flawed code is in the official GitHub main branch, git version *a39220048b3c9f2cca9a4d3a54604793c68eca7e*, and specifically in lines #125, #129, #133-136, #157, #221-225, #227-228, #231 of the file *autoattack/autoattack.py*. We will open source our updated code and encourage future stochastic defense methods to be evaluated against the fixed code. The code now can be found at: <https://anonymous.4open.science/r/auto-attack-595C/README.md>.

After the fix, AutoAttack’s accuracy dropped by up to 10 points, producing similar results to our PGD20-EoT10 test results. However, using AutoAttack on DPDDPM with $S = 1000$ took nearly 25 hours, five times longer than PGD20-EoT10, so we will use PGD20-EoT10 for the following test.

Table 13: AutoAttack (*Rand* version) and PGD20-EoT10 performance on DBP methods for CIFAR-10 (the lower the better). The original AutoAttack produces high accuracy, after fixing, it achieves similar results to PGD20+EoT10 attack.

Method	l_∞			l_2		
	AutoAttack	AutoAttack _{Ours}	PGD20-EoT10	AutoAttack	AutoAttack _{Ours}	PGD20-EoT10
DiffPure	62.11	56.25	55.96	81.84	76.37	75.78
DPDDPM-1000	57.81	46.88	48.63	71.68	71.09	72.27
DPDDIM-100	50.20	40.62	44.73	77.15	70.70	71.68

Table 14: Clean and robust accuracy on different DBP methods for CIFAR-10, evaluated with AutoAttack_{Ours} (*Rand* version). All methods show consistent improvement when fine-tuned with ADDT.

Method	Vanilla			Ours		
	Clean	l_∞	l_2	Clean	l_∞	l_2
DiffPure	89.26	56.25	76.37	89.94	58.20	77.34
DPDDPM	85.94	46.88	71.09	85.64	48.63	72.27
DPDDIM	88.38	40.62	70.70	88.77	44.73	71.68

647 **K Robustness under Deterministic White-box Setting**

648 We evaluate robustness of ADDT fine-tuned models under Deterministic White-box Setting. We
 649 assess robustness across varying NFE and compare the performance of vanilla models and ADDT fine-
 650 tuned models, with results showcased in Figure 9. ADDT consistently improves model performance
 651 at different NFE, particularly noticeable at lower NFE. This affirms that ADDT equips the diffusion
 model with the ability to directly counter adversarial perturbations.

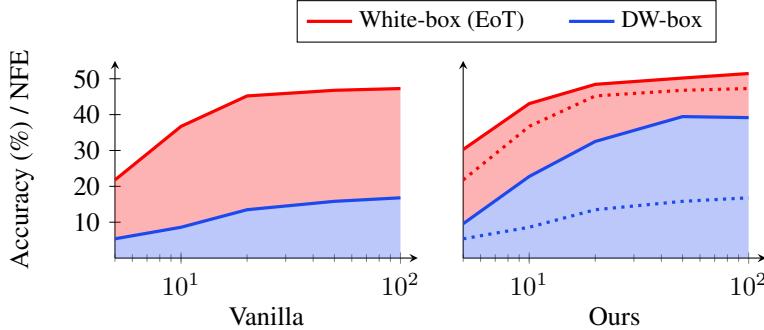


Figure 9: Revisiting Deterministic White-box Robustness. ADDT improves robustness under both White-box and Deterministic White-box setting, implying that ADDT strengthens the models’ ability to handle adversarial inputs. The dashed line on the right is the performance of the vanilla model.

652

653 **L Ablation Study of λ_{unit}**

654 In Section 6.2 we choose $\lambda_{unit}=0.03$ because most of the adversarial perturbations are in this range.
 655 We also provide an ablation study here, which shows that the performance of ADDT is insensitive to
 656 λ_{unit} and gets a consistent improvement.

Table 15: Ablation study of λ_{unit} , ADDT is insensitive to it and gets a consistent improvement

Attack type \NFE	λ_{unit}	50	100	200	500	1000
l_∞	Clean	21.78	36.72	45.21	46.78	47.27
	0.02	24.02	40.92	48.14	48.83	48.93
	0.03	30.27	43.07	48.44	50.20	51.46
	0.04	31.25	44.92	50.68	51.07	50.88
l_2	Clean	36.13	55.47	65.23	68.85	69.34
	0.02	41.99	61.72	67.48	69.82	70.31
	0.03	41.99	62.97	69.82	69.04	70.12
	0.04	49.02	64.45	69.24	69.53	69.92

657 **M Computational Cost of ADDT**

658 Optimizing DDPM and DDIM models through ADDT achieves near-optimal performance in 50
 659 epochs. This process takes only 12 hours on a four-GPU cluster of NVIDIA GeForce RTX 2080 Ti,
 660 matching the speed of traditional adversarial training and significantly faster than the latest adversarial
 661 training techniques that use diffusion models to augment the dataset [Wang et al., 2023].

662 A special feature of ADDT is its one-and-done training approach. After initial training, ADDT can
 663 protect various classifiers without the need for further fine-tuning, as shown in Table 4. This is
 664 different from adversarial classifier training, which requires individual training for each classifier.

665 **N Credibility of Our Paper**

666 The code was developed independently by two individuals and mutually verified, with consistent
 667 results achieved through independent training and testing. We will also make the code open-source
 668 and remain committed to advancing the field.

669 **O Broader Impact and Limitations**

670 Our work holds significant potential for positive societal impacts across various sectors, including
671 autonomous driving, facial recognition payment systems, and medical assistance. We are dedicated
672 to enhancing the safety and trustworthiness of global AI applications. However, there are potential
673 negative societal impacts, particularly concerning privacy protection, due to adversarial perturbations.
674 Nonetheless, we believe that the positive impacts generally outweigh the potential negatives. Regard-
675 ing the limitations, our approach could benefit from integrating insights from traditional adversarial
676 training methods [Zhang et al., 2019; Shafahi et al., 2019b; Wang et al., 2023], such as through more
677 extensive data augmentation and a refined ADDT loss design. Nevertheless, these limitations are
678 minor and do not significantly detract from the overall contributions of this paper.

679 **NeurIPS Paper Checklist**

680 **1. Claims**

681 Question: Do the main claims made in the abstract and introduction accurately reflect the
682 paper's contributions and scope?

683 Answer: [\[Yes\]](#)

684 Justification: The abstract and introduction provide a concise overview that aligns well with
685 the paper's contributions and delineates the scope effectively.

686 Guidelines:

- 687 • The answer NA means that the abstract and introduction do not include the claims
688 made in the paper.
- 689 • The abstract and/or introduction should clearly state the claims made, including the
690 contributions made in the paper and important assumptions and limitations. A No or
691 NA answer to this question will not be perceived well by the reviewers.
- 692 • The claims made should match theoretical and experimental results, and reflect how
693 much the results can be expected to generalize to other settings.
- 694 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
695 are not attained by the paper.

696 **2. Limitations**

697 Question: Does the paper discuss the limitations of the work performed by the authors?

698 Answer: [\[Yes\]](#)

699 Justification: The paper acknowledges the limitations of the authors' work, outlining areas
700 that require further research and improvement.

701 Guidelines:

- 702 • The answer NA means that the paper has no limitation while the answer No means that
703 the paper has limitations, but those are not discussed in the paper.
- 704 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 705 • The paper should point out any strong assumptions and how robust the results are to
706 violations of these assumptions (e.g., independence assumptions, noiseless settings,
707 model well-specification, asymptotic approximations only holding locally). The authors
708 should reflect on how these assumptions might be violated in practice and what the
709 implications would be.
- 710 • The authors should reflect on the scope of the claims made, e.g., if the approach was
711 only tested on a few datasets or with a few runs. In general, empirical results often
712 depend on implicit assumptions, which should be articulated.
- 713 • The authors should reflect on the factors that influence the performance of the approach.
714 For example, a facial recognition algorithm may perform poorly when image resolution
715 is low or images are taken in low lighting. Or a speech-to-text system might not be
716 used reliably to provide closed captions for online lectures because it fails to handle
717 technical jargon.
- 718 • The authors should discuss the computational efficiency of the proposed algorithms
719 and how they scale with dataset size.
- 720 • If applicable, the authors should discuss possible limitations of their approach to
721 address problems of privacy and fairness.
- 722 • While the authors might fear that complete honesty about limitations might be used by
723 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
724 limitations that aren't acknowledged in the paper. The authors should use their best
725 judgment and recognize that individual actions in favor of transparency play an impor-
726 tant role in developing norms that preserve the integrity of the community. Reviewers
727 will be specifically instructed to not penalize honesty concerning limitations.

728 **3. Theory Assumptions and Proofs**

729 Question: For each theoretical result, does the paper provide the full set of assumptions and
730 a complete (and correct) proof?

731 Answer: [NA]

732 Justification: The paper does not include theoretical results.

733 Guidelines:

- 734 • The answer NA means that the paper does not include theoretical results.
- 735 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 736 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 737 • The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- 738 • Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- 739 • Theorems and Lemmas that the proof relies upon should be properly referenced.

744 4. Experimental Result Reproducibility

745 Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

748 Answer: [Yes]

749 Justification: The paper provides comprehensive details necessary for replicating the primary experimental outcomes, supporting its main claims and conclusions effectively. This ensures transparency and facilitates further research in the field.

752 Guidelines:

- 753 • The answer NA means that the paper does not include experiments.
- 754 • If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- 755 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- 756 • Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- 757 • While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - 758 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - 759 (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - 760 (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - 761 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

784 5. Open access to data and code

785 Question: Does the paper provide open access to the data and code, with sufficient instruc-
786 tions to faithfully reproduce the main experimental results, as described in supplemental
787 material?

788 Answer: [Yes]

789 Justification: The paper includes comprehensive supplemental material with accessible
790 data and code, alongside detailed instructions, enabling accurate replication of the core
791 experimental findings.

792 Guidelines:

- 793 • The answer NA means that paper does not include experiments requiring code.
794
- 795 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
796
- 797 • While we encourage the release of code and data, we understand that this might not be
798 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
799 including code, unless this is central to the contribution (e.g., for a new open-source
benchmark).
- 800 • The instructions should contain the exact command and environment needed to run to
801 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
802
- 803 • The authors should provide instructions on data access and preparation, including how
804 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
805
- 806 • The authors should provide scripts to reproduce all experimental results for the new
807 proposed method and baselines. If only a subset of experiments are reproducible, they
808 should state which ones are omitted from the script and why.
809
- 810 • At submission time, to preserve anonymity, the authors should release anonymized
811 versions (if applicable).
812
- 813 • Providing as much information as possible in supplemental material (appended to the
paper) is recommended, but including URLs to data and code is permitted.

812 6. Experimental Setting/Details

813 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
814 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
815 results?

816 Answer: [Yes]

817 Justification: The paper provides comprehensive details on training and testing, including
818 data splits, hyperparameters and their selection process, ensuring a clear understanding of
819 the results.

820 Guidelines:

- 821 • The answer NA means that the paper does not include experiments.
822
- 823 • The experimental setting should be presented in the core of the paper to a level of detail
that is necessary to appreciate the results and make sense of them.
824
- 825 • The full details can be provided either with the code, in appendix, or as supplemental
material.

826 7. Experiment Statistical Significance

827 Question: Does the paper report error bars suitably and correctly defined or other appropriate
828 information about the statistical significance of the experiments?

829 Answer: [No]

830 Justification: The paper addresses the issue of experimental error by presenting the mean re-
831 sults from three separate trials. However, it does not provide error bars or other conventional
832 indicators of statistical variability.

833 Guidelines:

- 834 • The answer NA means that the paper does not include experiments.
835
- 836 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
dence intervals, or statistical significance tests, at least for the experiments that support
the main claims of the paper.
837

- 838 • The factors of variability that the error bars are capturing should be clearly stated (for
839 example, train/test split, initialization, random drawing of some parameter, or overall
840 run with given experimental conditions).
841 • The method for calculating the error bars should be explained (closed form formula,
842 call to a library function, bootstrap, etc.)
843 • The assumptions made should be given (e.g., Normally distributed errors).
844 • It should be clear whether the error bar is the standard deviation or the standard error
845 of the mean.
846 • It is OK to report 1-sigma error bars, but one should state it. The authors should
847 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
848 of Normality of errors is not verified.
849 • For asymmetric distributions, the authors should be careful not to show in tables or
850 figures symmetric error bars that would yield results that are out of range (e.g. negative
851 error rates).
852 • If error bars are reported in tables or plots, The authors should explain in the text how
853 they were calculated and reference the corresponding figures or tables in the text.

854 **8. Experiments Compute Resources**

855 Question: For each experiment, does the paper provide sufficient information on the com-
856 puter resources (type of compute workers, memory, time of execution) needed to reproduce
857 the experiments?

858 Answer: [Yes]

859 Justification: The paper does indeed offer comprehensive details on the computer resources
860 required, including the types of compute workers, memory specifications, and execution
861 time, ensuring the experiments are reproducible.

862 Guidelines:

- 863 • The answer NA means that the paper does not include experiments.
864 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
865 or cloud provider, including relevant memory and storage.
866 • The paper should provide the amount of compute required for each of the individual
867 experimental runs as well as estimate the total compute.
868 • The paper should disclose whether the full research project required more compute
869 than the experiments reported in the paper (e.g., preliminary or failed experiments that
870 didn't make it into the paper).

871 **9. Code Of Ethics**

872 Question: Does the research conducted in the paper conform, in every respect, with the
873 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

874 Answer: [Yes]

875 Justification: The research adheres to the ethical guidelines. It also addresses potential
876 societal impacts, aligning with the code's principles.

877 Guidelines:

- 878 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
879 • If the authors answer No, they should explain the special circumstances that require a
880 deviation from the Code of Ethics.
881 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
882 eration due to laws or regulations in their jurisdiction).

883 **10. Broader Impacts**

884 Question: Does the paper discuss both potential positive societal impacts and negative
885 societal impacts of the work performed?

886 Answer: [Yes]

887 Justification: The paper examines both the potential positive and negative societal impacts
888 of the work conducted.

889 Guidelines: The paper addresses both the potential positive and negative societal impacts of
890 the work.

- 891 • The answer NA means that there is no societal impact of the work performed.
- 892 • If the authors answer NA or No, they should explain why their work has no societal
893 impact or why the paper does not address societal impact.
- 894 • Examples of negative societal impacts include potential malicious or unintended uses
895 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
896 (e.g., deployment of technologies that could make decisions that unfairly impact specific
897 groups), privacy considerations, and security considerations.
- 898 • The conference expects that many papers will be foundational research and not tied
899 to particular applications, let alone deployments. However, if there is a direct path to
900 any negative applications, the authors should point it out. For example, it is legitimate
901 to point out that an improvement in the quality of generative models could be used to
902 generate deepfakes for disinformation. On the other hand, it is not needed to point out
903 that a generic algorithm for optimizing neural networks could enable people to train
904 models that generate Deepfakes faster.
- 905 • The authors should consider possible harms that could arise when the technology is
906 being used as intended and functioning correctly, harms that could arise when the
907 technology is being used as intended but gives incorrect results, and harms following
908 from (intentional or unintentional) misuse of the technology.
- 909 • If there are negative societal impacts, the authors could also discuss possible mitigation
910 strategies (e.g., gated release of models, providing defenses in addition to attacks,
911 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
912 feedback over time, improving the efficiency and accessibility of ML).

913 11. Safeguards

914 Question: Does the paper describe safeguards that have been put in place for responsible
915 release of data or models that have a high risk for misuse (e.g., pretrained language models,
916 image generators, or scraped datasets)?

917 Answer: [NA]

918 Justification: The paper does not involve the release of any high risk data or models.

919 Guidelines:

- 920 • The answer NA means that the paper poses no such risks.
- 921 • Released models that have a high risk for misuse or dual-use should be released with
922 necessary safeguards to allow for controlled use of the model, for example by requiring
923 that users adhere to usage guidelines or restrictions to access the model or implementing
924 safety filters.
- 925 • Datasets that have been scraped from the Internet could pose safety risks. The authors
926 should describe how they avoided releasing unsafe images.
- 927 • We recognize that providing effective safeguards is challenging, and many papers do
928 not require this, but we encourage authors to take this into account and make a best
929 faith effort.

930 12. Licenses for existing assets

931 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
932 the paper, properly credited and are the license and terms of use explicitly mentioned and
933 properly respected?

934 Answer: [Yes]

935 Justification: We have credited all original creators of the assets used in our paper, clearly
936 mentioning the asset versions, URLs, and specific licenses like CC-BY 4.0. For scraped
937 data, we complied with the source's copyright and terms of service.

938 Guidelines:

- 939 • The answer NA means that the paper does not use existing assets.
- 940 • The authors should cite the original paper that produced the code package or dataset.

- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No human subjects were involved in the research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer:[NA]

Justification: No human subjects were involved in the research.

Guidelines:

- 992 • The answer NA means that the paper does not involve crowdsourcing nor research with
993 human subjects.
994 • Depending on the country in which research is conducted, IRB approval (or equivalent)
995 may be required for any human subjects research. If you obtained IRB approval, you
996 should clearly state this in the paper.
997 • We recognize that the procedures for this may vary significantly between institutions
998 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
999 guidelines for their institution.
1000 • For initial submissions, do not include any information that would break anonymity (if
1001 applicable), such as the institution conducting the review.