
Threat Model-Agnostic Adversarial Defense using Diffusion Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep Neural Networks (DNNs) are highly sensitive to imperceptible malicious
2 perturbations, known as adversarial attacks. Following the discovery of this vulner-
3 ability in real-world imaging and vision applications, the associated safety concerns
4 have attracted vast research attention, and many defense techniques have been de-
5 veloped. Most of these defense methods rely on adversarial training (AT) – training
6 the classification network on images perturbed according to a specific threat model,
7 which defines the magnitude of the allowed modification. Although AT leads to
8 promising results, training on a specific threat model fails to generalize to other
9 types of perturbations. A different approach utilizes a preprocessing step to remove
10 the adversarial perturbation from the attacked image. In this work, we follow the
11 latter path and aim to develop a technique that leads to robust classifiers across
12 various realizations of threat models. To this end, we harness the recent advances
13 in stochastic generative modeling, and means to leverage these for sampling from
14 conditional distributions. Our defense relies on an addition of Gaussian i.i.d noise
15 to the attacked image, followed by a pretrained diffusion process – an architecture
16 that performs a stochastic iterative process over a denoising network, yielding
17 a high perceptual quality denoised outcome. The obtained robustness with this
18 stochastic preprocessing step is validated through extensive experiments on the
19 CIFAR-10 and CIFAR-10-C datasets, showing that our method outperforms the
20 leading defense methods under various threat models.

21 A Introduction

22 Deep neural network (DNN) image-classifiers are highly sensitive to malicious perturbations in which
23 the input image is slightly modified so as to change the classification prediction to a wrong class.
24 Amazingly, such attacks can be effective even with imperceptible changes to the input images. These
25 perturbations are known as adversarial attacks [10, 23, 37]. With the introduction of these DNN
26 classifiers to real-world applications, such as autonomous driving, this vulnerability has attracted vast
27 research attention, leading to the development of many attacks and robustification techniques.

28 Amongst the many types of adversarial attacks, the most common ones are norm-bounded to some
29 radius ϵ , where the norm L_p and the radius ϵ define a threat model. The attack is posed as an
30 optimization task in which one seeks the most effective deviation to the input image, δ , in terms of
31 modifying the classification output, while constraining this deviation to satisfy $\|\delta\|_p \leq \epsilon$. One way
32 to robustify a network against such attacks is by training it to correctly-classify attacked examples
33 from a specific threat model [25, 45, 11]. These methods, known as Adversarial Training (AT),
34 lead to state-of-the-art performance when trained and tested on the same threat model. However,
35 a well-known limitation of such methods is their poor generalization to unseen attacks, which is
36 discussed in length in [13, 2] as one of the unsolved problems of adversarial defense.

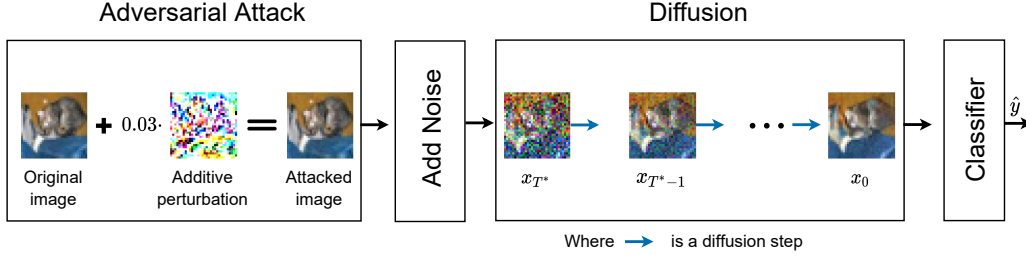


Figure 1: Our method flow. In the “Adversarial Attack” block, an attacker calculates the attack “Additive perturbation” and adds it to the “Original image” in order to create the “Attacked image”. As a preparation for the diffusion process, in the “Add Noise” block, we add an i.i.d Gaussian noise to the attacked image according to Equation 1. We proceed by feeding it into the “Diffusion” block, consisting of diffusion steps that include a denoising and an addition of a Gaussian noise. This effectively samples a new image from the diffusion model initialized by y , the noisy attacked image. Lastly, we feed the preprocessed obtained image to a classifier.

37 A different type of robustification techniques proposes a preprocessing step before feeding the image
 38 into the classifier [36, 30, 42, 12, 7, 15, 43]. Since an adversarial example can be seen as a summation
 39 of an image and an adversarial perturbation δ , using such a procedure to remove or even attenuate
 40 this second term is reasonable. The authors of [36, 30, 12, 7, 15, 43] use a generative model in the
 41 preprocessing phase in various ways. They either use the pretrained classifier directly or re-train a
 42 classifier on the generative model’s outputs. In general, these kind of methods are very appealing
 43 since they are capable of robustifying any publicly-available non-robust classifier and do not require
 44 a computational expensive specialized adversarial training. Furthermore, such methods are oblivious
 45 of the threat model being used.

46 In this work we introduce a novel and highly effective preprocessing robustification method for image
 47 classifiers. We choose a preprocessing-based approach based on a generative model since we aim to
 48 remove or weaken the adversarial perturbation while effectively projecting it onto the learned image
 49 manifold, where the classifier’s accuracy is likely to be high. While a generative model is typically
 50 used to sample from $p(x)$, the probability of images in general, our approach initializes this process
 51 with y at the appropriate diffusion step, where y is the noisy attacked image. This process effectively
 52 denoises the attacked image while targeting perfect perceptual quality [21, 28]. More specifically, we
 53 use a diffusion model - an iterative process that uses a pretrained MMSE (Minimum Mean Squared
 54 Error) denoiser and Langevin dynamics. The later involves an injection of Gaussian noise, which
 55 helps to robustify our samplers against attacks, even if they are aware of our defense strategy. Our
 56 method relies on a preprocessing model and a classification one, where both are trained independently
 57 on clean images. Hence, our architecture is inherently threat model agnostic, achieving robustness for
 58 unseen attacks. In our experiments we propose a way to evaluate the threat model-agnostic robustness
 59 by presenting two measurements. The first is the average on a wide range of attacks, and the second
 60 is the average across the unseen attacks. We consider the following threat models: $(L_\infty, \epsilon = 8/255)$,
 61 $(L_\infty, \epsilon = 16/255)$, $(L_2, \epsilon = 1)$, $(L_2, \epsilon = 2)$. In summary, our main contributions are:

- 62 • A novel stochastic diffusion-based preprocessing robustification is proposed, aiming to be a model-
 63 agnostic adversarial defense.
- 64 • The effectiveness of the proposed defense strategy is demonstrated in extensive experiments, showing
 65 state-of-the-art results.

66 B Our method

67 In this section we present our adversarial defense method, depicted in Figure 1. We start by adding
 68 noise to the attacked image, and then proceed by preprocessing the obtained image using a generative
 69 diffusion model, effectively projecting it onto the learned image manifold. The outcome of this
 70 diffusion is fed into a vanilla classifier, which is trained on the same image distribution that the
 71 diffusion model attempts to sample from. Thus, our framework is comprised of two main components
 72 – a denoiser that drives the diffusion model and a classifier.

73 Intuitively, we would like to sample images that are semantically close to an input image x by starting
 74 the diffusion process from some intermediate time step ($T^* < T$) rather than the beginning ($T^* = T$).
 75 Recall that x_T stands for a pure Gaussian noise, whereas x_{T^*} would be the noisy image we embark
 76 from. To this end, we modify the image to fit the diffusion model at this time step by applying
 77 Equation 1 – simply multiplying x by a scalar and adding an appropriate Gaussian noise, resulting
 78 in x_{T^*} . We feed this processed image into the diffusion model at time step T^* and complete the
 79 diffusion process, running with $t = T^*, T^* - 1, \dots, 0$, and outputting x_0 . Such a partial diffusion
 80 is similar to the image editing process presented in [26], and close in spirit to the posterior sampler
 81 that is discussed in [17]. We provide a comprehensive description of our method in Algorithm "Our
 82 preprocessing defense method" in the supplementary material.

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon; \quad \epsilon \sim \mathcal{N}(0, I). \quad (1)$$

83 An important hyperparameter for the success of our method is the initial diffusion depth T^* , since
 84 different values of it yield significant changes in x_0 . To better understand the importance of a careful
 85 choice of T^* , we intuitively analyze its effect. On the one hand, when starting from $T^* = T$,
 86 we sample a random image from the generative diffusion model, which obviously eliminates the
 87 adversarial perturbation. However, as the resulting image is independent of x , this will necessarily
 88 change class-related semantics of the image, which in turn would lead to misclassification. On
 89 the other hand, choosing $T^* = 0$ results in the same input image x , which does not remove the
 90 perturbation from the image, hence probably leading to misclassification as well. In other words, we
 91 need to choose T^* that balances the trade-off between cleaning the adversarial noise, and keeping
 92 the semantic properties of the input image x . Choosing such T^* that successfully balances these
 93 properties is crucial to the success of our adversarial defense algorithm.

94 We utilize the above described sampling algorithm with one goal in mind – sampling an image that is
 95 not contaminated with an adversarial attack while keeping it semantically similar to the original input
 96 image x . We believe that our algorithm is suited for this task because the Gaussian noise injections are
 97 much larger than the adversarial perturbation. Hence, the noise overshadows the adversarial attack,
 98 reducing its effect. This leads to a sampling process that answers both of our demands, removal of
 99 the contamination while remaining semantically close to x .

100 As mentioned previously, our method is comprised of a diffusion model denoiser and a classifier,
 101 both trained on clean images. This framework is very useful from a practical point of view, since we
 102 can utilize publicly available pretrained models to a completely different task than they were trained
 103 on – adversarial defense. The fact that these models were trained without adversarial attacks in mind
 104 gives our method a significant advantage – it is inherently threat model-agnostic. This essentially
 105 avoids the challenged generalization to unseen attacks problem [13, 2], according to which classifiers
 106 trained on a specific adversarial threat model are vulnerable to attacks under a different threat regime.

107 C Experiments

108 We proceed by empirically demonstrating the improved performance attained by our proposed
 109 adversarial defense method. We compare our method to various state-of-the-art (SoTA) methods on
 110 white box attacks. Additional experiments are reported in the supplementary material.

111 Throughout our experiments, we use the pretrained diffusion model from [34] and a vanilla
 112 classifier, both trained on clean images from CIFAR-10 [22] train set (50,000 examples). More
 113 specifically, we set the diffusion model maximal depth to $T^* = 140$ and the sub-sequence of
 114 the time steps to $\tau = \{T^*, T^* - 10, \dots, 10, 0\}$. In addition, we use a WideResNet-28-10 [44]
 115 architecture as our classifier and evaluate the performance on the CIFAR-10 test set (10,000 examples).
 116

117 In order to evaluate our defense in the best way, we must use white-box attacks. This allows us
 118 to estimate the worst examples for our defense. But, conducting white box experiments requires
 119 differentiating through all of the diffusion time steps and also through the classifier. Since it requires
 120 infeasible computation power we use a standard PGD attack - identical to the algorithm presented in
 121 [25]¹. Moreover, in order to have a fair comparison, we use the exact same attack to evaluate all of
 122 the different defenses.

¹<https://github.com/MadryLab/robustness>

Table 1: CIFAR-10 robust accuracies under white + EOT attacks. For every compared method, we state the threat model that was used in training in the first column Trained Threat Model (TTM) column. The next four columns are the four different threat models used for evaluation. The next two columns are the two averages that we use for evaluation, Average without Training (AwT), and Average of All (AoA). In the last column we state the classifier architecture that is used.

Method	TTM	Attack				AwT	AoA	Architecture
		L_∞		L_2				
		8/255	16/255	1	2			
AT [25]	$L_\infty, \epsilon = 8/255$	54.23	19.20	32.34	04.99	18.84	27.69	rn-50
	$L_2, \epsilon = 0.5$	34.25	02.99	41.55	05.72	21.13	21.13	rn-50
Trades [45]	$L_\infty, \epsilon = 8/255$	55.79	23.18	32.51	05.01	20.23	29.12	wrn-34-10
Gowal et al. [11]	$L_\infty, \epsilon = 8/255$	66.35	34.81	41.87	09.62	28.77	38.16	wrn-28-10
	$L_2, \epsilon = 0.5$	47.08	13.12	52.71	14.85	31.94	31.94	wrn-70-16
PAT - [24]		44.07	22.33	46.65	23.33	34.01	34.01	rn-50
Ours		51.05	37.76	50.75	19.23	39.70	39.70	wrn-28-10

123 **C.1 CIFAR-10 experimets**

124 Next, we compare our method to baseline state-of-the-art (SoTA) methods, under PGD attacks using
 125 four different threat models – ($L_2, \epsilon = 1$), ($L_2, \epsilon = 2$), ($L_\infty, \epsilon = 8/255$), ($L_\infty, \epsilon = 16/255$)- more
 126 details are given in supplementary material. To assess the generalization ability to unseen attacks,
 127 we average the results in two ways: (i) *Average of All*: accuracy average of all the attacks; and (ii)
 128 *Average of Unseen Attack*: accuracy average of the attacks not seen at training time (if applicable).
 129 While the first is a simple average that also considers the performance on the attack used in training
 130 time, the second showcases the generalization capabilities to unseen attacks. Note that because our
 131 method is not trained on any threat model, (i) and (ii) are the same. As can be seen in Table 1,
 132 adversarial training methods excel on the specific threat model that they trained on. However, they
 133 generalize poorly, as discussed in [2, 13], while our method achieves SoTA performance in both of
 134 the examined metrics.

135 **D Conclusion**

136 This work presents a novel preprocessing defense mechanism against adversarial attacks, based on a
 137 generative diffusion model. Since this generative model relies on pretraining on clean images, it has
 138 the capability to generalize to unseen attacks. We evaluate our method across different attacks and
 139 demonstrate its superior performance. Our method can be used to defend against any attack, and does
 140 not require retraining the vanilla classifier.

References

- [1] A. Athalye, N. Carlini, and D. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pages 274–283. PMLR, 2018.
- [2] T. Bai, J. Luo, J. Zhao, B. Wen, and Q. Wang. Recent advances in adversarial training for adversarial robustness. *arXiv preprint arXiv:2102.01356*, 2021.
- [3] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow. Thermometer encoding: One hot way to resist adversarial examples. In *International Conference on Learning Representations*, 2018.
- [4] N. Carlini and D. Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.
- [5] J. Cohen, E. Rosenfeld, and Z. Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019.
- [6] F. Croce, S. Gowal, T. Brunner, E. Shelhamer, M. Hein, and T. Cemgil. Evaluating the adversarial robustness of adaptive test-time defenses. *arXiv preprint arXiv:2202.13711*, 2022.
- [7] Y. Du and I. Mordatch. Implicit generation and modeling with energy based models. *Advances in Neural Information Processing Systems*, 32, 2019.
- [8] L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, B. Tran, and A. Madry. Adversarial robustness as a prior for learned representations. *arXiv preprint arXiv:1906.00945*, 2019.
- [9] R. Ganz and M. Elad. Bigroc: Boosting image generation via a robust classifier, 2021.
- [10] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [11] S. Gowal, C. Qin, J. Uesato, T. Mann, and P. Kohli. Uncovering the limits of adversarial training against norm-bounded adversarial examples. *arXiv preprint arXiv:2010.03593*, 2020.
- [12] W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, M. Norouzi, and K. Swersky. Your classifier is secretly an energy based model and you should treat it like one. *arXiv preprint arXiv:1912.03263*, 2019.
- [13] D. Hendrycks, N. Carlini, J. Schulman, and J. Steinhardt. Unsolved problems in ml safety. *arXiv preprint arXiv:2109.13916*, 2021.
- [14] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- [15] M. Hill, J. Mitchell, and S.-C. Zhu. Stochastic security: Adversarial defense using long-run dynamics of energy-based models. *arXiv preprint arXiv:2005.13525*, 2020.
- [16] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [17] Z. Kadkhodaie and E. P. Simoncelli. Solving linear inverse problems using the prior implicit in a denoiser. *arXiv preprint arXiv:2007.13640*, 2020.
- [18] D. Kang, Y. Sun, D. Hendrycks, T. Brown, and J. Steinhardt. Testing robustness against unforeseen adversaries. *arXiv preprint arXiv:1908.08016*, 2019.
- [19] B. Kawar, M. Elad, S. Ermon, and J. Song. Denoising diffusion restoration models. *arXiv preprint arXiv:2201.11793*, 2022.
- [20] B. Kawar, G. Vaksman, and M. Elad. Snips: Solving noisy inverse problems stochastically. *Advances in Neural Information Processing Systems*, 34, 2021.

- 184 [21] B. Kawar, G. Vaksman, and M. Elad. Stochastic image denoising by sampling from the posterior
185 distribution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
186 pages 1866–1875, 2021.
- 187 [22] A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 188 [23] A. Kurakin, I. Goodfellow, and S. Bengio. Adversarial machine learning at scale. *arXiv preprint*
189 *arXiv:1611.01236*, 2016.
- 190 [24] C. Laidlaw, S. Singla, and S. Feizi. Perceptual adversarial robustness: Defense against unseen
191 threat models. *arXiv preprint arXiv:2006.12655*, 2020.
- 192 [25] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models
193 resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- 194 [26] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon. SDEdit: Guided image
195 synthesis and editing with stochastic differential equations. In *International Conference on*
196 *Learning Representations*, 2021.
- 197 [27] A. Q. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *International*
198 *Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- 199 [28] G. Ohayon, T. Adrai, G. Vaksman, M. Elad, and P. Milanfar. High perceptual quality image
200 denoising with a posterior sampling cgan. In *Proceedings of the IEEE/CVF International*
201 *Conference on Computer Vision*, pages 1805–1813, 2021.
- 202 [29] E. Raff, J. Sylvester, S. Forsyth, and M. McLean. Barrage of random transforms for adversarially
203 robust defense. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
204 *Recognition*, pages 6528–6537, 2019.
- 205 [30] P. Samangouei, M. Kabkab, and R. Chellappa. Defense-gan: Protecting classifiers against
206 adversarial attacks using generative models. *arXiv preprint arXiv:1805.06605*, 2018.
- 207 [31] S. Santurkar, A. Ilyas, D. Tsipras, L. Engstrom, B. Tran, and A. Madry. Image synthesis with a
208 single (robust) classifier. *Advances in Neural Information Processing Systems*, 32, 2019.
- 209 [32] A. Shamir, O. Melamed, and O. BenShmuel. The dimpled manifold model of adversarial
210 examples in machine learning. *arXiv preprint arXiv:2106.10151*, 2021.
- 211 [33] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning
212 using nonequilibrium thermodynamics. In *International Conference on Machine Learning*,
213 pages 2256–2265. PMLR, 2015.
- 214 [34] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint*
215 *arXiv:2010.02502*, 2020.
- 216 [35] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution.
217 *Advances in Neural Information Processing Systems*, 32, 2019.
- 218 [36] Y. Song, T. Kim, S. Nowozin, S. Ermon, and N. Kushman. Pixeldefend: Leveraging generative
219 models to understand and defend against adversarial examples. *arXiv preprint arXiv:1710.10766*,
220 2017.
- 221 [37] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus.
222 Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- 223 [38] F. Tramer, N. Carlini, W. Brendel, and A. Madry. On adaptive attacks to adversarial example
224 defenses. *Advances in Neural Information Processing Systems*, 33:1633–1645, 2020.
- 225 [39] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness may be at odds
226 with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- 227 [40] C. Xie, J. Wang, Z. Zhang, Z. Ren, and A. Yuille. Mitigating adversarial effects through
228 randomization. *arXiv preprint arXiv:1711.01991*, 2017.

- 229 [41] W. Xu, D. Evans, and Y. Qi. Feature squeezing: Detecting adversarial examples in deep neural
230 networks. *arXiv preprint arXiv:1704.01155*, 2017.
- 231 [42] Y. Yang, G. Zhang, D. Katabi, and Z. Xu. Me-net: Towards effective adversarial robustness
232 with matrix estimation. *arXiv preprint arXiv:1905.11971*, 2019.
- 233 [43] J. Yoon, S. J. Hwang, and J. Lee. Adversarial purification with score-based generative models.
234 In *International Conference on Machine Learning*, pages 12062–12072. PMLR, 2021.
- 235 [44] S. Zagoruyko and N. Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*,
236 2016.
- 237 [45] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. El Ghaoui, and M. Jordan. Theoretically principled
238 trade-off between robustness and accuracy. In *International conference on machine learning*,
239 pages 7472–7482. PMLR, 2019.

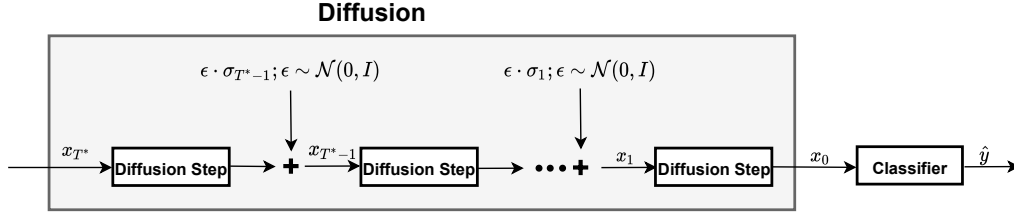


Figure 2: Our method incorporates a diffusion model and a classifier. In every diffusion step, we add Gaussian noise multiplied by the corresponding σ_t , which is a user-controlled hyperparameter. The variables x_{T^*}, \dots, x_1 constitute the MCMC, and the last step’s output of the diffusion model x_0 , is the final output, to be sent to the classifier.

240 E Background

241 E.1 Adversarial robustness

242 Since the discovery of the phenomenon of adversarial examples in neural networks [10, 23, 37],
 243 classifiers’ robustness has been extensively studied. Numerous works have been focusing on new
 244 methods for constructing adversarial examples and/or defending from them. In the following we bring
 245 the very fundamental results referring to adversarial defense and attack methods, as a background to
 246 our work.

247 Let us start with how adversarial attacks are created. Given an image x and a classifier $f(\cdot)$, an
 248 adversarial attack is a small norm-bounded perturbation δ , added to the input image x , that leads to
 249 its misclassification. There exist several mainstream settings for crafting adversarial examples that
 250 differ from each other in their assumptions regarding the defense method’s characteristics and the
 251 access to the model and its gradients. We describe below such key attack configurations.

252 *White-Box Attacks* are applied when the attacker has full access to the full system architecture
 253 (including both the classifier and the defense mechanism), which is assumed to be differentiable.
 254 This is a rich and a widely used group of attacks that contains some of the most common ones, such
 255 as Fast Gradient Signed Method (FGSM) [10], Projected Gradient Decent (PGD)[25] and CW [4].
 256 While there exist numerous white-box attack strategies, PGD is the cornerstone of their most modern
 257 embodiments. It is an iterative gradient-based algorithm that increases the classifier’s loss in each
 258 step by perturbing the input data. We describe PGD in Algorithm 1 below.

Algorithm 1 L_∞ -based Projected Gradient Descent

Input classifier $f(\cdot)$, input x , target label y , norm radius ϵ , step size α , number of steps N

```

1: procedure PGD
2:    $\delta \leftarrow 0$ 
3:   for  $i$  in  $1 : N$  do
4:      $\delta \leftarrow \Pi_\epsilon(\delta + \alpha \cdot \text{sign}(\nabla_x \text{Loss}(f(x + \delta), y)))$ 
5:   end for
6: end procedure

```

259 The operator Π_ϵ is a projection onto the L_p norm of radius ϵ . In the L_∞ case, Π_ϵ is just the clamp
 260 operation into $[-\epsilon, \epsilon]$.

261 Since white-box attacks have assumptions that do not always hold, they can not be used in every setup.
 262 For example, such a setup can be a defense method that relies on a non-differentiable preprocessing.
 263 Since white-box attacks are gradient-based, they are likely to fail in this case. Another example is
 264 stochastic preprocessing, which poses a challenging configuration for white-box attacks. This stems
 265 from the fact that the ideal crafted attack might not be optimal during inference due to randomness.
 266 In order to better adjust gradient-based adversarial attacks to such scenarios, alternative approaches
 267 were developed, as we describe hereafter.

268 *Grey-Box Attack* is used when the attacker has access to the classifier but not to the preprocessing
 269 model defending it, $g(\cdot)$. This approach is limited due to the fact that the attack in such a case

270 is constructed upon $f(\cdot)$ while being evaluated with $f(g(\cdot))$. As a consequence, the malicious
 271 perturbation created is necessarily sub-optimal and thus less effective.

272 *Backward Pass Differentiable Approximation (BPDA) Attack* [1] is an attack method for cases in
 273 which the preprocessing function $g(\cdot)$ is non-differentiable or impractical to differentiate, implying
 274 that $f(g(\cdot))$ is not differentiable as well. In many cases we can invoke the assumption that $g(x) \approx x$,
 275 reflecting the fact that preprocessing methods do not perform significant modifications to the input
 276 images, but rather try to remove the already small malicious perturbations. In order to attack such
 277 architecture we use the forward pass of the preprocessing $g(\cdot)$ and approximate its derivative with
 278 I , producing $\nabla_x f(g(x)) \approx \nabla_{g(x)} f(g(x))$. With this in place, the attacker can perform white-box
 279 attacks without completely disregarding the preprocessing steps.

280 *Expectation-Over-Transformation (EOT) Attack* [1] is used when the preprocessing step $g(\cdot)$ is
 281 stochastic. Attacking such a method is harder for gradient-based methods, since the crafted deviation
 282 vector δ might not remain optimal during inference due to the randomness. EOT calculates the
 283 attack’s gradients by $\nabla_x \mathbb{E}[f(g(x))] = \mathbb{E}[\nabla_x f(g(x))]$, differentiating through both the classifier and
 284 preprocessing with an expectation. In practice, EOT empirically approximates the expectation with a
 285 fixed number of drawn samples from $g(x)$.

286 We move now to discuss adversarial defense approaches. In the past few years, numerous such
 287 methods were proposed to improve the robustness of classifiers to adversarial attacks. While there
 288 are many types of robustification algorithms, we focus below on two such families.

289 *Adversarial Training (AT) Defense* proposes to utilize adversarial examples during the training process
 290 of the classifier. More specifically, the idea is to train the model to classify such examples correctly.
 291 Several recent works [25, 45, 11] follow this line of reasoning, leading to the current state-of-the-art
 292 in robustifying classifiers.

293 *Preprocessing* is a substantially different type of robustification method that relies on a preceding
 294 operation on the classifier’s input as its name suggests. Since adversarial examples contain small
 295 imperceptible perturbations, using preprocessing steps to “clean” them seems to be an intuitive
 296 step. Many works rely on various generative models for such preprocessing [36, 30, 7, 15, 43]. More
 297 specifically, these models are used to project the attacked image into a valid clean one in its vicinity,
 298 with the hope that the processed image is more likely to be classified correctly.

299 E.2 Diffusion models

300 Diffusion models [33, 16, 35] are Markov Chain Monte Carlo (MCMC)-based generative techniques,
 301 which consist of a chain of images x_0, x_1, \dots, x_T of the same size as the given image x . These
 302 methods are based on two closely related processes. The first is the forward process of gradually
 303 adding Gaussian noise to the data according to a decaying variance schedule parametrized by
 304 $1 > \alpha_0 > \alpha_1 > \dots > \alpha_T > 0$. The following defines this chain of steps, for $t = 1, 2, \dots, T$ where
 305 x_0 is the given clean image x :

$$q(x_t|x_{t-1}) := \mathcal{N}\left(\sqrt{\frac{\alpha_t}{\alpha_{t-1}}}x_{t-1}, \left(1 - \frac{\alpha_t}{\alpha_{t-1}}\right)I\right) \quad (2)$$

306 Posed differently, the forward process can be described as a simple weighting between the image x_0
 307 and a Gaussian noise vector,

$$q(x_t|x_0) = \mathcal{N}(\sqrt{\alpha_t}x_0, (1 - \alpha_t)I), \quad (3)$$

308 so we can express x_t as

$$x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon; \quad \epsilon \sim \mathcal{N}(0, I). \quad (4)$$

309 When α_t is close to zero, x_t is close to a pure standard Gaussian noise, independent of x_0 . Thus, we
 310 can set $x_T \sim \mathcal{N}(0, I)$ as initialization for the backward process, which is explained next.

311 The second and the more intricate process is the backward direction, which gradually removes the
 312 noise from the image. Intuitively, this stage denoises the image by peeling layers of noise gradually. A

313 key ingredient in this process is a pretrained noise estimator neural network, $\epsilon_\theta(x_t, t)$. This denoiser
 314 serves as an approximation to the score function $\nabla \log p(x)$ [17], bringing the knowledge about the
 315 image statistics into this sampling procedure. The noise estimator is conditioned on the time t , trying
 316 to estimate the noise ϵ of the latent variable x_t . Sampling, or generating an image, is performed by
 317 iteratively applying the following update rule for $t = T, T - 1, \dots, 0$:

$$x_{t-1} = \sqrt{\alpha_{t-1}} \left(\frac{x_t - \sqrt{1 - \alpha_t} \epsilon_\theta(x_t, t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \epsilon_\theta(x_t, t) + \sigma_t \epsilon_t \quad (5)$$

318 where the first term is a denoising stage – an estimation of x_0 , while the second term stands for an
 319 attenuated version of the estimated additive noise in x_t . $\sigma_t \epsilon_t$ is a stochastic addition, where σ_t is a
 320 hyperparameter controlling the stochasticity of the process, and $\epsilon_t \sim \mathcal{N}(0, I)$.

321 The sampling process posed in Equation (5) tends to be very slow, requiring T (≈ 1000) passes
 322 through the denoising network. Methods for speeding up this process are discussed in [27, 34, 19].
 323 There are various use-cases for diffusion models beyond image synthesis. The ones relevant to our
 324 work are discussed in [26, 21, 20, 19] where inverse problems are being considered. Following [26],
 325 instead of sampling from the ideal image distribution $p(x)$, the diffusion process we implement is
 326 initialized with x_{T^*} , where x_{T^*} ² is the given noisy image. Thus, the outcome x_0 can be considered
 327 as a stochastic high perceptual quality denoising of x_{T^*} .

Algorithm 2 Our preprocessing defense method

Input image x , maximum depth T^* , diffusion model denoiser $\epsilon_\theta(\cdot, \cdot)$,
 variance schedule $[\alpha_T, \dots, \alpha_0]$, stochasticity hyperparameters $[\sigma_T, \dots, \sigma_1]$,

- 1: **procedure** SAMPLING
- 2: $\epsilon_{T^*} \sim \mathcal{N}(0, I)$
- 3: $x_{T^*} \leftarrow \sqrt{\alpha_{T^*}} x + \sqrt{1 - \alpha_{T^*}} \epsilon_{T^*}$
- 4: **for** t **in** $[T^*, T^* - 1, \dots, 1]$ **do**
- 5: $\tilde{x}_{t-1} \leftarrow \frac{x_t - \sqrt{1 - \alpha_t} \epsilon_\theta(x_t, t)}{\sqrt{\alpha_t}}$
- 6: $\epsilon_t \sim \mathcal{N}(0, I)$
- 7: $x_{t-1} \leftarrow \sqrt{\alpha_{t-1}} \tilde{x}_{t-1} + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \epsilon_\theta(x_t, t) + \sigma_t \epsilon_t$
- 8: **end for**
- 9: **return** x_0
- 10: **end procedure**

328 F Experiments

329 F.1 Synthetic Dataset Experiments

330 We create a synthetic 2D dataset (see Figure 3) and investigate the effect of a diffusion process on the
 331 decision boundaries of the classification. The dataset consists of two classes – red and blue points –
 332 consisting altogether of 10,000 examples, drawn from two mixtures of Gaussians, each consisting of
 333 4 concentrated groups. We train a fully connected neural network model to classify this data, having
 334 10 layers of width 128. The training is done via 5,000 epochs. As for the diffusion preprocess, we
 335 use an analytic score-function $\nabla \log p(x)$ of the known distribution, following the work of [35]. We
 336 set $T^* = 10$ and values of α in the range $[0.1, 1]$.

337 After training the classifier, we calculate its decision rule and present it in Figure 3a, where the
 338 background colors represent the predicted label. As can be seen, the classifier achieves perfect
 339 performance, as all the red points are located in the red zone, and all the blue ones are surrounded by a
 340 blue background. Nevertheless, the classifier decision boundaries are very close to the data, which is
 341 a well-known phenomenon of vanilla classifiers [32]. This illustrates why small perturbations to the
 342 data, such as adversarial attacks, can change the classification decision from the correct to the wrong
 343 ones.

²More on the relation between T and T^* is given below.



Figure 3: Decision boundary comparison between a vanilla classifier with and without our method on a 2D synthetic dataset.

Table 2: CIFAR-10 robust accuracies of preprocessing methods under the following attack: grey-box, BPDA + EOT, white-box PGD. All using the same threat model $L_\infty, \epsilon = 8/255$.

Defense	Attack	Base Classifier		Preprocessed	
		Clean	Adversarial	Clean	Adversarial
ADP [43]	grey-box	95.60	00.00	86.39	80.49
Ours	grey-box	95.60	00.00	86.28	82.33
ADP [43]	BPDA+EOT	95.60	00.00	86.39	44.79
Ours	BPDA+EOT	95.60	00.00	86.28	77.65
ADP [43]	white-box	95.60	00.00	86.39	31.42
Ours	white-box	95.60	00.00	86.28	63.40

344 When applying our preprocessing scheme, our method leads to a larger margin between the data
 345 points and the decision boundaries, as can be seen in Figure 3b. These results are encouraging
 346 because in the adversarial attack regime, every data point is allowed to be perturbed with an ϵ norm ball
 347 around it. When the decision boundaries are far enough from the data points, an ϵ -bounded attack
 348 would necessarily fail.

349 F.2 CIFAR-10 experiments

350 First, we compare our method to ADP [43], a leading preprocessing method, using the following
 351 attacks: grey-box, BPDA+EOT, and white-box, where the EOT is approximated over 20 repetitions.
 352 As can be seen in Table 2, our method outperforms ADP by up to 32.86%. We should note that the
 353 results are lower than presented in [43], this was also observed in [6].

354 When deploying the proposed diffusion defense, two critical parameters should be discussed - the
 355 choice of T^* (referred to as depth) and the time-step skips to use. In this Subsection we discuss the
 356 effect of both.

357 We start by showing the influence of the depth of the diffusion model on the robust accuracy. As we
 358 change the maximal depth of the diffusion model T^* , we depict the robust accuracy obtained by our
 359 method, and present it in Figure 4. As discussed in Section Our Method, the diffusion depth controls
 360 the trade-off between clearing the attack perturbation and sampling an image that is semantically
 361 similar to the input image x . We track the diffusion model behavior as we increase the diffusion
 362 model’s first step. When setting T^* to a shallow diffusion step, we effectively sample images that
 363 are closer to the input image x , and since the image is contaminated by a malicious attack, the
 364 classification accuracy is low. As we increase the depth we reach a sweet-spot in which we clean the
 365 malicious perturbation while keeping a small perceptual distance to x , which leads to the highest
 366 accuracy. When the depth is too big, we clear the attack but lose perceptual similarity to x , and the
 367 accuracy is reaching 10%, meaning that we sample random images.

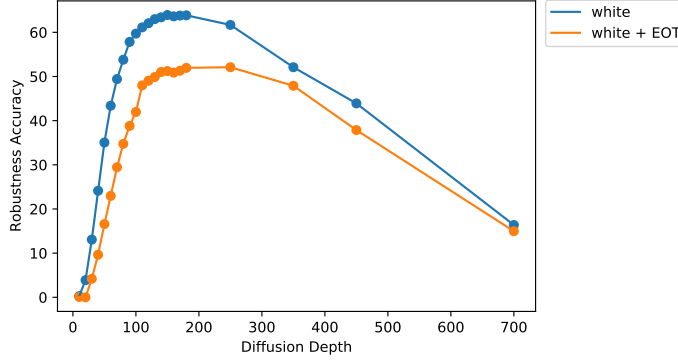


Figure 4: The obtained robust accuracy under white box attacks as a function of the max depth T^* of the diffusion model. There are two graphs, both are attacked using the same threat model $L_\infty, \epsilon = 8/255$, the first is the robust accuracy under white-box attack, and the other refers to a white-box + EOT.

Table 3: CIFAR-10 robust accuracies under white + EOT attacks. We present two samplings of the diffusion model time steps. The first uses $\tau = \{T^*, T^* - 10, \dots, 10, 0\}$ while the second applies a full sampling $\tau = \{T^*, T^* - 1, \dots, 1, 0\}$. We compare the two sampling performance. In the ‘‘Attack’’ columns we present the accuracy under different threat models. The last two columns are two averages used for evaluation: Average without Training (AwT), and Average of All (AoA). It was evaluated on the first 1000 test images of CIFAR10

Method	Attack				AwT	AoA
	L_∞		L_2			
	8/255	16/255	1	2		
Ours	61.54	43.66	63.64	43.56	53.10	53.10
Ours - full sampling	64.14	44.06	63.74	47.15	54.77	54.77

368 We now move to explore the influence of the skips to the time-steps in the diffusion process. Attacking
 369 our preprocessing method necessarily consumes a lot of time and memory, making it hard to break, as
 370 indeed claimed in [15]. This is due to the fact that an attack process requires keeping a computational
 371 graph of all the time steps of the diffusion process for computing derivatives. In contrast, our defense
 372 mechanism is lighter, as no derivatives are required, and only T^* forward passes through the denoiser
 373 are performed.

374 When evaluating our defense method under the strongest known attack, white-box + EOT, we must
 375 lighten further our protection by reducing the number of diffusion steps. This is done by using
 376 only 1/10 of the DDIM diffusion steps [34], requiring all-together 14 steps. For uniformity of our
 377 experiments, we use this sub-sequence of steps for all attacks.

378 We should note that if the proposed preprocess diffusion is applied in full (no subsampling), this
 379 would increase both the attack and defense runtime and memory consumptions by a factor of 10.
 380 Such an approach would not worsen the robust accuracy, and perhaps even improve it, as can be seen
 381 in Table 3. Both these effects have one clear conclusion – when using our defense in practice, we can
 382 increase the diffusion model sampling, harming the attacker, while preserving the robust accuracy.

383 Attack structure

384 Working with adversarial perturbation of images has the advantages of enabling the analysis of
 385 the attack δ , better understanding it, and getting an intuition about it. When an attack changes the
 386 classification prediction of an image, one might expect the perceptual structure of the image to change

387 accordingly, just as is accomplished in order to change a human’s prediction. However, this is not
 388 always the case when fooling a deep-neural-network classifier.

389 A geometrical explanation for this phenomenon is given in [32], showing that trained vanilla classifiers
 390 tend to produce decision boundaries that are nearly parallel to the data manifold. As such, fooling the
 391 network amounts to a very small step orthogonal to this manifold, thus having no “visual meaning”.
 392 In contrast, robust classifiers behave differently, exhibiting Perceptual Aligned Gradients (PAG)
 393 [39, 31, 8, 9].

394 White-box attacks of the form we consider in this work are based on computing the gradients of the
 395 attacked classifier. Therefore, when a classifier exhibits a PAG property in its gradients, this would
 396 imply a highly desired robustness behavior. Armed with this insight, we consider the following
 397 question: Given a system comprising of both the vanilla classifier and our diffusion-based defense
 398 mechanism, does this overall system have PAG?

399 We answer the above question and present some empirical evidence of this phenomenon in Figure
 400 5. In the first row we show several original images from CIFAR-10. In the second row we present
 401 a white-box attack on a vanilla classifier, an attack lacking perceptual meaning. In the third row
 402 we present white-box + EOT attack under our method, exhibiting PAG - the obtained gradients
 403 concentrate on the object, aiming to modify its appearance. When attacking the defended classifier,
 404 the attacker use white-box + EOT, an attack that was crafted for stochastic defenses. Every attack’s
 405 step is the expectation over multiple realizations of the defense.

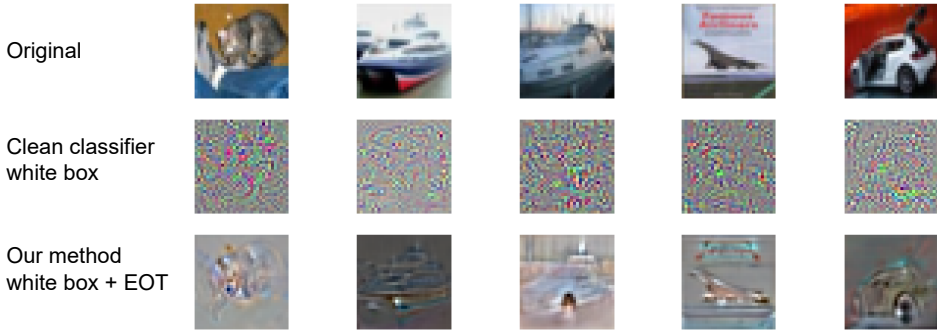


Figure 5: The attack δ structure of white-box+EOT attack, L_2 norm, radius $\epsilon = 1$. First row: Five CIFAR-10 images. Second row: The attack δ under a white-box attack, where the attacked classifier is a vanilla one. Third row: The attack on our method, where we preprocess the image before inputting into a vanilla classifier.

406 **Robustness to CIFAR-10-C perturbations**

407 In most of our discussion we focused on a robustness to norm- bounded attacks. We turn now
 408 to introduce a robust classification under attacks that are based on augmentation. These refer to
 409 modifications of the image in various ways such as motion blur, zoom blur, snow, JPEG compression,
 410 contrast variation, etc. CIFAR-10-C [14] is such a corrupted images dataset that was created by
 411 performing numerous augmentations on CIFAR-10 [22] dataset. CIFAR-10-C is commonly used for
 412 evaluating the robustness performance under broad attacks.

413 As our method is inherently attack agnostic, it is natural to evaluate it on this class of attacks.
 414 We compare our method versus other leading techniques, achieving state-of-the-art results. This
 415 experiment requires adjustment of the diffusion model maximal depth parameter T^* . When we set
 416 $T^* \in [30, 90]$, we outperform the other methods, as depicted in Figure 6.

417 **Computational resources**

418 Our proposed defense method relies on an application of a diffusion model as a preprocessing stage for
 419 purifying adversarial perturbations. To perform a gradient-based attack, one needs to backpropagate

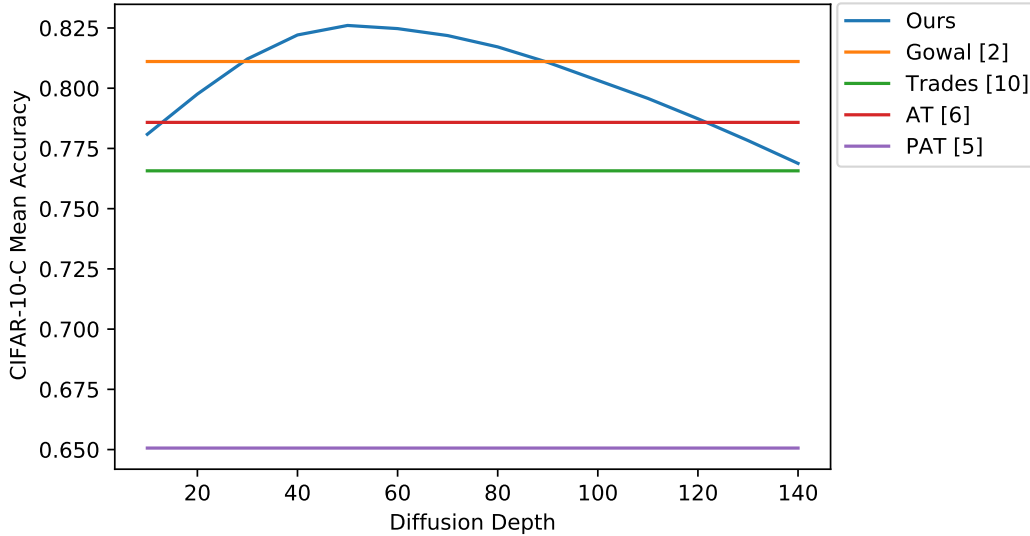


Figure 6: Robustness accuracy under CIFAR-10-C as a function of the diffusion model maximal depth T^* . We compare our method with the results reported in [11, 45, 25, 24].

420 the gradients through the classifier and the diffusion model. This process is very expensive, both in
 421 terms of memory and computations, since the attacker needs to keep the entire computational graph
 422 in memory and backpropagate from the classifier through all of the diffusion time steps.

423 When evaluating our defense method under our most challenging attack, white-box + EOT, we must
 424 further lighten our approach by reducing the number of diffusion steps. We do so by using only 1/10
 425 of the diffusion steps, i.e., 14 times instead of 140. This reduction decreases the computational needs
 426 and enables us to perform such an attack, using 8 NVIDIA A4000 GPUs. As shown in Table 3 the
 427 robust accuracy of our method is slightly reduced, while significantly improving the computational
 428 cost and achieving state-of-the-art performance.

429 G Related work

430 The goal of preprocessing methods is to clean the adversarial attacks from the input images, leading to
 431 correct prediction by deep neural network classifier. Preliminary work on preprocessing defense meth-
 432 ods include rescaling [40], thermometer encoding [3], feature squeezing [41], GAN for reconstruction
 433 [30], ensemble of transformations [29], addition of Gaussian noise [5] and mask and reconstruction
 434 [42]. It was shown by [1, 38] that such preprocessing, even if it includes stochasticity and non-
 435 differentiability, can be broken when evaluated properly by adjusting the projected-gradient-descent
 436 attack, using backward-pass-differentiable-approximation and expectation-over-transformation algo-
 437 rithms. A new preprocessing group of work has recently emerged, trying to utilize Energy-Based-
 438 Model (EBM) to the task of cleaning adversarial perturbation from images. The intuition is that
 439 generative models are capable of sampling images from the image manifold, hopefully projecting
 440 attacked images that were deviated from the image manifold, back onto it. To this end, some EBM
 441 preprocessing methods were developed: purification by pixelCNN [36], restore corrupt image with
 442 EBM [7] and density aware classifier [12]. Most recent methods includes: long-run Langevin sam-
 443 pling [15] and gradient ascent score based-model [43]. In contrast to many of these methods that
 444 require retraining the classifier, our method does not have this requirement, the diffusion model and
 445 classifier are both pretrained on clean images.

446 Defense to unseen attacks methods: Recently, an attention for defense to unseen attacks has emerged.
 447 Previous methods that include Adversarial Training (AT) do not generalize well to unseen attacks,
 448 as shown in [13, 2]. For this end, a new robustness evaluation metric to unseen attacks was suggested
 449 [18]. Moreover, the authors of [24] suggested perceptual-adversarial-training, which takes into
 450 account the perceptual similarity, leading to a new method that generalizes to unseen attacks.