
On the Limitations of Stochastic Pre-processing Defenses

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Defending against adversarial examples remains an open problem. A common
2 belief is that randomness at inference increases the cost of finding adversarial
3 inputs. An example of such a defense is to apply a random transformation to
4 inputs prior to feeding them to the model. In this paper, we empirically and
5 theoretically investigate such stochastic pre-processing defenses and demonstrate
6 that they are flawed. First, we show that most stochastic defenses are weaker than
7 previously thought; they lack sufficient randomness to withstand even standard
8 attacks like projected gradient descent. This casts doubt on a long-held assumption
9 that stochastic defenses invalidate attacks designed to evade deterministic defenses
10 and force attackers to integrate the Expectation over Transformation (EOT) concept.
11 Second, we show that stochastic defenses confront a trade-off between adversarial
12 robustness and model invariance; they become less effective as the defended model
13 acquires more invariance to their randomization. Future work will need to decouple
14 these two effects. Our code is available in the supplementary material.

15 1 Introduction

16 Machine learning models are vulnerable to adversarial examples [4, 29], where an adversary can
17 add imperceptible perturbations to the input of a model and change its prediction [5, 22]. Their
18 discovery has motivated a wide variety of defense approaches [6, 8, 10, 25, 34, 35] along with the
19 evaluation of their adversarial robustness [2, 24, 32]. Current evaluations mostly rely on adaptive
20 attacks [2, 32], which require significant modeling and computational efforts. However, even when
21 the attack succeeds, such evaluations may not always reveal the fundamental weaknesses of an
22 examined defense. Without awareness of the underlying weaknesses, subsequent defenses may still
23 conduct inadvertently weak adaptive attacks; this leads to overestimated robustness.

24 One popular class of defenses that demonstrates the above is the stochastic pre-processing defense,
25 which relies on applying randomized transformations to inputs to provide robustness [10, 35]. Despite
26 existing attack techniques designed to handle randomness [2, 3], there is an increasing effort to
27 improve these defenses through a larger randomization space or more complicated transformations.
28 For example, BaRT [24] employs 25 transformations, where the parameters of each transformation
29 are further randomized. Due to the complexity of this defense, it was only broken recently (three
30 years later) by Sitawarin et al. [28] with a complicated adaptive attack. Still, it is unclear how future
31 defenses can avoid the pitfalls of existing defenses, largely because these pitfalls remain unknown.

32 In this paper, we investigate stochastic pre-processing defenses and explain their limitations both
33 empirically and theoretically. First, we revisit previous stochastic pre-processing defenses and explain
34 why such defenses are broken. We show that most stochastic defenses are not sufficiently randomized
35 to invalidate standard attacks designed for deterministic defenses. Second, we study recent stochastic
36 defenses that exhibit more randomness and show that they also face key limitations. In particular,

37 we identify a trade-off between their robustness and the model’s invariance to their transformations.
38 These defenses achieve a notion of robustness that results from reducing the model’s invariance to
39 the applied transformations. We outline our findings below. These findings suggest future work to
40 find new ways of using randomness that decouples these two effects.

41 **Most stochastic defenses lack sufficient randomness.** Although Athalye et al. [2] and Tramèr
42 et al. [32] have demonstrated the ineffectiveness of several stochastic defenses with techniques like
43 Expectation over Transformation (EOT) [3], it remains unclear whether and why EOT is required
44 to break them. A commonly accepted explanation is that EOT computes the “correct gradients” of
45 models with randomized components [2, 32], yet the necessity of such correct gradients has not been
46 explicitly discussed. To fill this gap, we examine a long-held assumption that stochastic defenses
47 invalidate standard attacks designed for deterministic defenses.

48 Specifically, we revisit stochastic pre-processing defenses previously broken by EOT and examine
49 their robustness *without* applying EOT. Interestingly, we find that most stochastic defenses lack
50 sufficient randomness to withstand even standard attacks (that do not integrate any strategy to capture
51 model randomness) like projected gradient descent (PGD) [22]. We then conduct a systematic
52 evaluation to show that applying EOT is only beneficial when the defense is sufficiently randomized.
53 Otherwise, standard attacks already perform well and the randomization’s robustness is overestimated.

54 **Trade-off between adversarial robustness and model invariance.** When stochastic pre-processing
55 defenses do have sufficient randomness, they must fine-tune the model using augmented training data
56 to preserve utility in the face of randomness added. We characterize this procedure by the model’s
57 *invariance* to the applied defense, where we identify a trade-off between the model’s robustness
58 (provided by the defense) and its invariance to the applied defense. Stochastic pre-processing defenses
59 become less effective when their defended model acquires more invariance to their transformations.

60 On the theoretical front, we present a theoretical setting where this trade-off provably exists. We show
61 from this trade-off that stochastic pre-processing defenses provide robustness by inducing variance on
62 the defended model, and must take back such variance to recover utility. We verify this trade-off with
63 empirical evaluations on realistic datasets, models, and defenses. We observe that robustness drops
64 when the defended model is fine-tuned on data processed by its defense to acquire higher invariance.

65 2 Related Work

66 **Stochastic Defenses.** Defending against adversarial examples remains an open problem, where a
67 common belief is that inference-time randomness increases the cost of finding adversarial inputs. Early
68 examples of such stochastic defenses include input transformations [10] and rescaling [35]. These
69 defenses were broken by Athalye et al. [2] using techniques like EOT [3] to capture randomness. After
70 that, more stochastic defenses were proposed but with inadvertently weak evaluations [23, 25, 34, 36],
71 which were found ineffective by Tramèr et al. [32]. Subsequent stochastic defenses resort to larger
72 randomization space like BaRT [24], which was only broken recently by Sitawarin et al. [28].
73 Randomized smoothing [6, 16, 17] leverages randomness to certify robustness. In this work, instead of
74 designing adaptive attacks for individual defenses, we focus on the general stochastic pre-processing
75 defenses and demonstrate their limitations.

76 **Trade-offs for Adversarial Robustness.** The trade-offs associated with adversarial robustness have
77 been widely discussed in the literature. Prior work identified trade-offs between robustness and
78 accuracy [33, 37] and even between the robustness against different types of adversarial examples [12,
79 31]. Recent work also investigated the trade-off between the model’s robustness and invariance to
80 input transformations, such as circular shifts [27] and rotations [13]. These trade-offs characterize
81 a standalone model’s own property – the model itself is less robust to adversarial examples when
82 it becomes more invariant to certain transformations, without any defense. Our setting, however,
83 is orthogonal to such analysis; the model we consider is protected by a stochastic pre-processing
84 defense, and we aim to characterize the performance of the added defense to the model.

85 3 Preliminaries

86 **Notation.** Let $f : \mathcal{X} \rightarrow \mathbb{R}^C$ denote the classifier with pre-softmax outputs, where $\mathcal{X} = [0, 1]^d$ is
87 the input space with d dimensions and C is the number of classes. We then consider a stochastic

88 pre-processing defense $t_\theta : \mathcal{X} \rightarrow \mathcal{X}$, where θ is the random variable drawn from some randomization
 89 space Θ that parameterizes the defense. The defended classifier can be written as $f_\theta(\mathbf{x}) := f(t_\theta(\mathbf{x}))$.

90 Let $F(\mathbf{x}) := \arg \max_{i \in \mathcal{Y}} f_i(\mathbf{x})$ denote the classifier that returns the predicted label, where f_i is the
 91 output of the i -th class and $\mathcal{Y} = [C]$ is the label space. Similarly, we use F_θ and $f_{\theta,i}$ to denote the
 92 prediction and class-output of the stochastic classifier f_θ . Since this classifier returns varied outputs
 93 for a fixed input, it determines the final prediction by aggregating n independent inferences with
 94 strategies like majority vote. We discuss these strategies and the choice of n in Appendix A.1.

95 **Adversarial Examples.** Given an image $\mathbf{x} \in \mathcal{X}$ and a classifier F , the adversarial example
 96 $\mathbf{x}' := \mathbf{x} + \delta$ is visually similar to \mathbf{x} but either misclassified (i.e., $F(\mathbf{x}') \neq F(\mathbf{x})$) or classified as
 97 a target class y' chosen by the attacker (i.e., $F(\mathbf{x}') = y'$). Attack algorithms generate adversarial
 98 examples by searching for δ such that \mathbf{x}' fools the classifier while minimizing δ under some distance
 99 metrics; for instance, the ℓ_p norm constraint $\|\delta\|_p \leq \epsilon$ for a perturbation budget ϵ .

100 **Projected Gradient Descent (PGD).** PGD [22] is one of the most established attacks to evaluate
 101 adversarial example defenses. Given a benign example \mathbf{x}^0 and its ground-truth label y , each iteration
 102 of the untargeted PGD attack (with ℓ_∞ norm budget ϵ) can be formulated as

$$\mathbf{x}^{i+1} \leftarrow \mathbf{x}^i + \alpha \cdot \text{sgn}\{\nabla \mathcal{L}(f_\theta(\mathbf{x}^i), y)\}, \quad (1)$$

103 where α is the step size, \mathcal{L} is the loss function, and each iteration is projected to the ℓ_∞ ball around
 104 \mathbf{x}^0 of radius ϵ . We use PGD- k to denote the PGD attack with k steps. We outline formulations for
 105 other settings and norms in Appendix A.2.

106 **Expectation over Transformation (EOT).** Since the classifier f_θ is stochastic, the defense evalua-
 107 tion literature [2, 32] argues that attacks should target the *expectation* of the gradient using Expectation
 108 over Transformation (EOT) [3], which reformulates the PGD attack as

$$\mathbf{x}^{i+1} \leftarrow \mathbf{x}^i + \alpha \cdot \text{sgn}\left\{\mathbb{E}_{\theta \sim \Theta} \left[\nabla \mathcal{L}(f_\theta(\mathbf{x}^i), y)\right]\right\} \approx \mathbf{x}^i + \alpha \cdot \text{sgn}\left\{\frac{1}{m} \sum_{j=1}^m \nabla \mathcal{L}(f_{\theta_j}(\mathbf{x}^i), y)\right\}, \quad (2)$$

109 where m is the number of samples to estimate the expectation and $\theta_j \stackrel{\text{iid}}{\sim} \Theta$ are sampled parameters
 110 for the defense. We use EOT- m to denote the EOT technique with m samples at each PGD step.

111 In addition, for a fair comparison among attacks with different PGD steps and EOT samples, we
 112 quantify the attack’s strength by its total number of gradient computations. For example, attacks
 113 using PGD- k and EOT- m will have strength $k \times m$. Although white-box attacks are typically not
 114 constrained in this way, it allows for a fair comparison when attacks have finite computing resources
 115 (e.g., when EOT is not parallelizable). We discuss more about this quantification in Appendix A.3.

116 4 Most Stochastic Defenses Lack Sufficient Randomness

117 Athalye et al. [2] and Tramèr et al. [32] demonstrate adaptive evaluation of stochastic defenses with
 118 the application of EOT. However, it remains unclear why EOT is required to break these stochastic
 119 defenses. While a commonly accepted explanation is that EOT computes the “correct gradients” of
 120 models with randomized components [2, 32], the necessity of such correct gradients has not been
 121 explicitly discussed. To fill this gap, we revisit stochastic defenses previously broken by EOT and
 122 examine their robustness *without* applying EOT. Interestingly, we find that applying EOT is mostly
 123 *unnecessary* when evaluating existing stochastic defenses.

124 **Case Study: Random Rotation.** We start with a simple
 125 stochastic defense that randomly rotates the input image for
 126 $\theta \in [-90, 90]$ degrees (chosen at uniform) before classifica-
 127 tion. This defense is representative for most pre-processing
 128 defenses [10, 24, 35]. We evaluate this defense on 1,000 Im-
 129 ageNet images with PGD- k and EOT- m under the constraint
 130 $k \times m = 50$, as discussed in Section 3. All attacks use max-
 131 imum ℓ_∞ perturbation $\epsilon = 8/255$ with step size chosen from
 132 $\alpha \in \{1/255, 2/255\}$. The results are shown in Table 1, where

Table 1: Evaluation of the random rotation with PGD- k and EOT- m .

Attacks	k	m	Success Rate
Untargeted	10	5	100%
	50	1	100%
Targeted	10	5	99.0%
	50	1	99.0%

Table 2: The missing ablation study of adaptive evaluations of stochastic defenses in the literature. Notations: attack iterations k , EOT samples m , learning rate α , number of gradient queries $k \times m$. The details of these defenses and their evaluation settings are in Appendix B.

Defenses	Original Adaptive Evaluation (w/ EOT)					Our Ablation Study (w/o EOT)				
	k	m	α	$k \times m$	Success Rate	k	m	α	$k \times m$	Success Rate
Guo et al. [10]	1,000	30	0.1	30,000	100%	1,000	1	0.001	1,000	99.0%
Xie et al. [35]	1,000	30	0.1	30,000	100%	200	1	0.1	200	100%
Dhillon et al. [8]	500	10	0.1	5,000	100%	500	1	0.1	500	100%
Xiao et al. [34]	100	1,000	0.01	100,000	100%	40,000	1	0.1/255	40,000	98.4%
Roth et al. [25]	100	40	0.2/255	4,000	100%	4,000	1	0.1/255	4,000	96.1%

133 PGD-50 performs equally well as PGD-10 combined with EOT-5. This observation suggests that
 134 *some stochastic defenses are already breakable without applying EOT*, casting doubt on a long-held
 135 assumption that stochastic defenses simply invalidate attacks designed for deterministic defenses.

136 **Comprehensive Evaluations.** We then extend the above case study to other stochastic defenses
 137 evaluated in the literature. Specifically, we replicate the (untargeted) adaptive evaluation of stochastic
 138 defenses from Athalye et al. [2] and Tramèr et al. [32] with their official implementation. We only
 139 change the attack’s hyper-parameters (e.g., number of iterations and learning rate) and disable EOT
 140 by setting its number of samples to one ($m = 1$), which avoids potential implementation flaws if
 141 removed from the source code. The comparison between evaluations with and without applying EOT
 142 is summarized in Table 2, which serves as a missing ablation study of adaptive evaluations in the
 143 literature. The experimental settings are identical within each row (detailed in Appendix B).

144 Interestingly, we find it *unnecessary* to break these defenses with EOT, as long as the standard
 145 attack runs for more iterations with a smaller learning rate. For such defenses, standard iterative
 146 attacks already contain an *implicit expectation* across iterations to capture the limited randomness.
 147 This observation implies that most stochastic defenses lack sufficient randomness to withstand even
 148 standard attacks designed for deterministic defenses. Therefore, increasing randomness becomes a
 149 promising approach to enhancing stochastic defenses, as adopted by recent defenses [6, 24]. Note
 150 that this ablation study only aims to inspire potential ways of enhancing stochastic defenses; it does
 151 not invalidate EOT for stronger adaptive evaluations of stochastic defenses.

152 5 Trade-offs between Robustness and Invariance

153 When stochastic pre-processing defenses *do have* sufficient randomness, they must ensure that the
 154 utility of the defended model is preserved in the face of randomness. To achieve high utility, existing
 155 defenses mostly rely on augmentation invariance through *trained invariance* [21]. In such a case, the
 156 invariance is achieved by applying the defense’s randomness to the training data so as to guide the
 157 model in learning their transformations. For defenses based on stochastic pre-processor t_θ , each data
 158 sample from the dataset gets augmented with t_θ sampled from the randomization space Θ , and the
 159 risk is minimized over such augmented data.

160 The defended classifier $F_\theta(x) := F(t_\theta(x))$ is invariant under the randomization space Θ if

$$161 F(t_\theta(x)) = F(x), \quad \forall \theta \in \Theta, x \in \mathcal{X}. \quad (3)$$

162 As we can observe from the definition, invariance has direct implications on the performance of
 163 stochastic pre-processing defenses. If the classifier is invariant under the defense’s randomization
 164 space Θ as is defined in Equation (3), then the defense should not work – computing the model and
 165 its gradients over randomization $\theta \in \Theta$ is the same as if t_θ was not applied at all. This observation
 166 suggests a direct coupling between invariance and performance of the defense: the more invariant,
 167 hence performant, the model is under a given randomization space, the less protection such a defense
 168 would provide. In this section, we present a simple theoretical setting where this coupling provably
 exists; detailed computations and illustrations of this setting are in Appendix C.1.

169 **Binary Classification Task.** We consider a class-balanced dataset \mathcal{D} consisting of input-label pairs
 170 (x, y) with $y \in \{-1, +1\}$ and $x|y \sim \mathcal{N}(y, 1)$, where $\mathcal{N}(\mu, \sigma^2)$ is a normal distribution with mean μ
 171 and variance σ^2 . Moreover, an ℓ_∞ -bounded adversary perturbs the input with a small ϵ to fool the
 172 classifier. We quantify the classifier’s robustness by its robust accuracy, i.e., the ratio of correctly
 173 classified samples that remain correct after being perturbed by the adversary. Note that while we use
 174 $x + \epsilon$ to denote the perturbed input, its actual value can take within the range $[x - \epsilon, x + \epsilon]$.

175 *Undefended Classification.* We start with the optimal linear classifier $F(x) := \text{sgn}(x)$ without any
 176 defense. This classifier attains robust accuracy

$$\Pr[F(x + \epsilon) = y \mid F(x) = y] = \frac{\Pr[F(x + \epsilon) = y \wedge F(x) = y]}{\Pr[F(x) = y]} = \frac{\Phi(1 - \epsilon)}{\Phi(1)}, \quad (4)$$

177 where Φ is the cumulative distribution function of $\mathcal{N}(0, 1)$.

178 *Defended Classification.* We then try to improve adversarial robustness by introducing a stochastic
 179 pre-processing defense $t_\theta(x) := x + \theta$, where $\theta \sim \mathcal{N}(1, \sigma^2)$ is the random variable parameterizing
 180 the defense. This defense characterizes common pre-processing defenses that enforce randomness
 181 while shifting the input distribution. Here, the processed input follows a shifted distribution $t_\theta(x) \sim$
 182 $\mathcal{N}(y + 1, 1 + \sigma^2)$. As a result, the defended classifier $F_\theta(x) = \text{sgn}(x + \theta)$ has robust accuracy

$$\Pr[F_\theta(x + \epsilon) = y \mid F_\theta(x) = y] = \frac{\Pr[F_\theta(x + \epsilon) = y \wedge F_\theta(x) = y]}{\Pr[F_\theta(x) = y]} = \frac{\Phi'(-\epsilon) + \Phi'(2 - \epsilon)}{\Phi'(0) + \Phi'(2)}, \quad (5)$$

183 where $\Phi'(x) := \Phi(x/\sqrt{1 + \sigma^2})$ is the cumulative distribution function of $\mathcal{N}(0, 1 + \sigma^2)$. At this
 184 point, we have not fit the classifier on processed inputs. Due to its lack of invariance, the defended
 185 classifier has low utility yet higher robust accuracy than the undefended one in Equation (4).

186 *Defended Classification (Trained Invariance).* As discussed above, one critical step of stochastic
 187 pre-processing defenses is to preserve the defended model’s utility by minimizing the risk over
 188 augmented data $t_\theta(x)$, which leads to a new defended classifier $F_\theta^+(x) = \text{sgn}(x + \theta - 1)$. As a
 189 result, this new defended classifier achieves higher invariance with robust accuracy

$$\Pr[F_\theta^+(x + \epsilon) = y \mid F_\theta^+(x) = y] = \frac{\Pr[F_\theta^+(x + \epsilon) = y \wedge F_\theta^+(x) = y]}{\Pr[F_\theta^+(x) = y]} = \frac{\Phi'(1 - \epsilon)}{\Phi'(1)}, \quad (6)$$

190 which is less robust than the previous less-invariant classifier F_θ in Equation (5). However, one may
 191 observe that this classifier, though loses some robustness compared with F_θ , is still more robust
 192 than the original undefended classifier F in Equation (4). This part of robustness comes from the
 193 changed data distribution due to the defense’s randomness. It shows that we have not achieved perfect
 194 invariance to the defense’s randomness, thus gaining some robustness at the cost of utility.

195 *Defended Classification (Perfect Invariance).* Furthermore, these defenses usually leverage majority
 196 vote to obtain stable predictions, which finally produces a perfectly invariant defended classifier

$$F_\theta^*(x) = \text{sgn}\left\{\frac{1}{n} \sum_{i=1}^n F_{\theta_i}^+(x)\right\} = \text{sgn}\left\{\frac{1}{n} \sum_{i=1}^n \text{sgn}(x + \theta_i - 1)\right\} \rightarrow \text{sgn}(x) = F(x), \quad (7)$$

197 where $\theta_i \stackrel{\text{iid}}{\sim} \mathcal{N}(1, \sigma^2)$ are sampled parameters. In such a case, the defended classifier reduces to the
 198 original undefended classifier with the original robust accuracy:

$$\Pr[F_\theta^*(x + \epsilon) = y \mid F_\theta^*(x) = y] = \Pr[F(x + \epsilon) = y \mid F(x) = y] = \frac{\Phi(1 - \epsilon)}{\Phi(1)}. \quad (8)$$

199 **Summary.** The above theoretical setting illustrates how stochastic pre-processing defenses first induce
 200 variance on the binary classifier we consider to provide adversarial robustness in Equation (5), and
 201 how they finally take back such variance in Equations (6) and (8) to recover utility. In Appendix C.2,
 202 we extend this coupling to a general trade-off where we can control the invariance. This trade-off
 203 demonstrates that stochastic pre-processing defenses provide robustness by explicitly reducing the
 204 model’s invariance to added randomized transformations.

205 6 Experiments

206 Our experiments are designed to answer the following two questions.

207 Q1: What properties make applying EOT beneficial when evaluating stochastic defenses?

208 We show that applying EOT is only beneficial when the defense is sufficiently randomized; otherwise
 209 standard attacks already perform well and leave no room for EOT to improve.

210 **Q2: What is the limitation of stochastic defenses when they do have sufficient randomness?**

211 We show a trade-off between the stochastic defense’s robustness and the model’s invariance to
212 the defense itself. Such defenses become less effective when the defended model achieves higher
213 invariance to their randomness, as required to preserve utility under the defense.

214 6.1 Experimental Settings

215 **Datasets & Models.** We conduct all experiments on ImageNet [26] and ImageNette [9]. For Image-
216 Net, our test data consists of 1,000 images randomly sampled from the validation set. ImageNette is
217 a ten-class subset of ImageNet, and we test on its validation set. We adopt various ResNet [11] models.
218 For defenses with low randomness, we evaluate them on ImageNet with a pre-trained ResNet-50 with
219 Top-1 accuracy 75.9%. For defenses with higher randomness (thus requiring fine-tuning), we switch
220 to ImageNette and a pre-trained ResNet-34 with Top-1 accuracy 96.9% to reduce the training cost like
221 previous work [28]. These models are fine-tuned on the training data processed by tested defenses.
222 As a special case, we also evaluate randomized smoothing on ImageNet using the ResNet-50 models
223 from Cohen et al. [6]. More details of datasets and models can be found in Appendices D.1 and D.2.

224 **Defenses & Metrics.** We focus on stochastic defenses allowing us to increase randomness: random-
225 ized smoothing [6] and BaRT [24]. For randomized smoothing, we vary the variance of the added
226 Gaussian noise. For BaRT, we vary the number κ of applied randomized transformations. Note that
227 we have evaluated other stochastic defenses and discussed their low randomness in Section 4. We
228 measure the defense’s performance by the defended model’s *benign accuracy* and the attack’s *success*
229 *rate*, all evaluated with majority vote over $n = 500$ predictions. The attack’s success rate is the ratio
230 of samples that do not satisfy the attack’s objective prior to the attack but satisfy it after the attack.
231 For example, we discard samples that were misclassified before being perturbed in untargeted attacks.
232 Details of the evaluated defenses can be found in Appendix D.3.

233 **Attacks.** We evaluate defenses with standard PGD combined with EOT and focus on the ℓ_∞ -bounded
234 adversary with a perturbation budget $\epsilon = 8/255$ in both untargeted and targeted settings. We only
235 conduct adaptive evaluations, where the defense is included in the attack loop with non-differentiable
236 components captured by BPDA [2]. We also utilize AutoPGD [7] to avoid selecting the step size
237 when it is computationally expensive to repeat some experiments. More details of the attack’s setting
238 and implementation can be found in Appendix D.4.

239 6.2 Q1: Evaluate the Benefits of Applying EOT under Different Settings

240 In Section 4, we showed that standard attacks are sufficient to break most stochastic defenses due to
241 their lack of randomness. Here, we aim to understand what properties make applying EOT beneficial
242 when evaluating stochastic defenses. We design a systematic evaluation of stochastic defenses with
243 different level of randomness and check if applying EOT improves the attack.

244 **Stochastic Defenses with Low Randomness.** We start with BaRT’s noise injection defense, which
245 perturbs the input image with noise of distributions and parameters chosen at random. While this
246 defense has low randomness, it yields meaningful results. We evaluate this defense with various
247 combinations of PGD and EOT. The performance of untargeted and targeted attacks is shown in
248 Figure 1. We test multiple step sizes and summarize their best results (discussed in Appendix E.1).

249 In this case, standard PGD attacks are already good enough when the defense has insufficient
250 randomness, leaving no space for improvements from EOT. In Figure 1f, both (1) PGD-10 combined
251 with EOT-10 and (2) PGD-100 combined with EOT-1 have near 100% success rates. This result is
252 consistent with our observations in Section 4 in both untargeted and targeted settings¹.

253 **Stochastic Defenses with Higher Randomness.** We then examine the randomized smoothing
254 defense that adds Gaussian noise to the input image. Although this defense was originally proposed
255 for certifiable adversarial robustness, we adopt it to evaluate how randomness affects the benefits of
256 applying EOT. Similarly, we evaluate this defense with PGD and EOT of different settings with a
257 focus on the *targeted* attack. The results are shown in Figure 2.

¹The only caveat is that targeted attacks are more likely to benefit from EOT, as their objectives are stricter and may have better performance with gradients of higher precision.

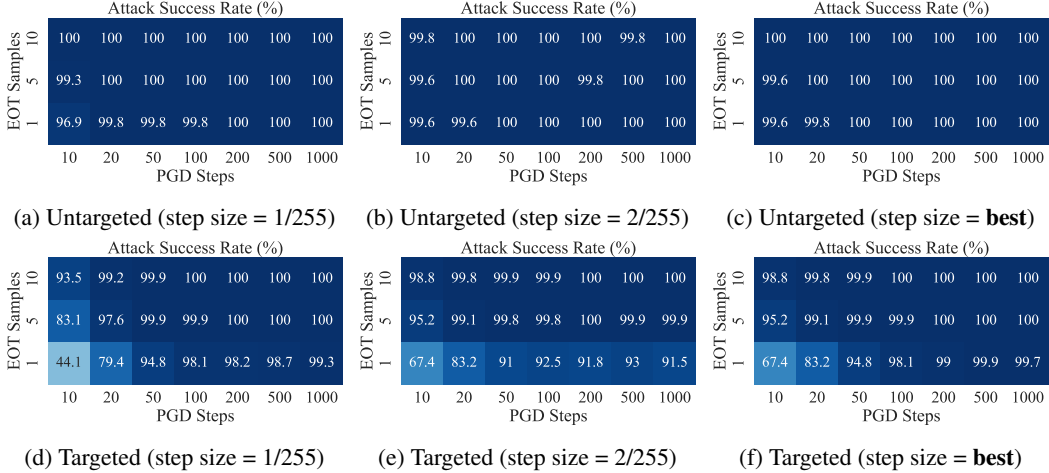


Figure 1: Evaluation of BaRT’s noise injection defense on ImageNet. Standard PGD without applying EOT (i.e., applying EOT-1) is already good enough, leaving limited space for EOT to improve.

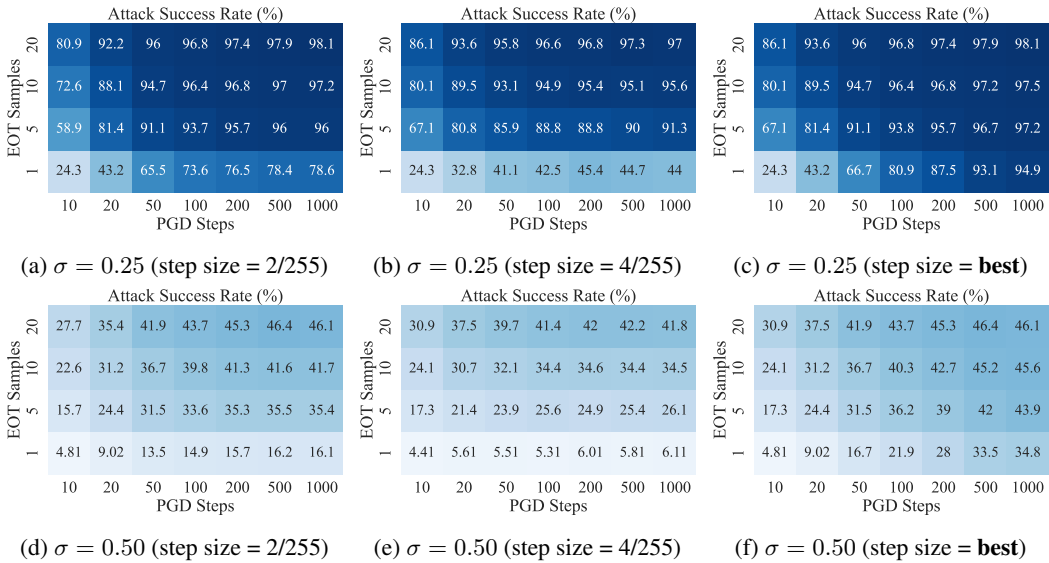


Figure 2: Evaluation of randomized smoothing on ImageNet (targeted attacks). PGD performs well on lower variance ($\sigma = 0.25$) if running for more steps. For a larger variance ($\sigma = 0.50$), applying EOT starts to improve the attack significantly (for a fixed number of gradient computations).

258 We observe that EOT starts to improve the attack when the defense has a higher level of randomness.
 259 For a fixed number of PGD steps, applying EOT significantly improves the attack in most of the
 260 settings. For a fixed attack strength (i.e., number of gradient computations), applying EOT always
 261 outperforms standalone PGD. In Figure 2f, for example, PGD-100 combined with EOT-10 is 5.5%
 262 higher than PGD-1,000 with EOT-1 (40.3% vs. 34.8%).

263 **Takeaways.** Applying EOT is only beneficial when the defense has sufficient randomness, such as
 264 randomized smoothing with $\sigma = 0.5$. This observation suggests that stochastic defenses only make
 265 standard attacks suboptimal when they have sufficient randomness. However, most existing stochastic
 266 defenses did not achieve this criterion, as we showed in Section 4.

267 6.3 Q2: Evaluate the Trade-off between Robustness and Invariance

268 In Section 5, we present a theoretical setting where the trade-off between robustness and invariance
 269 provably exists; stochastic defenses become less robust when the defended model achieves higher

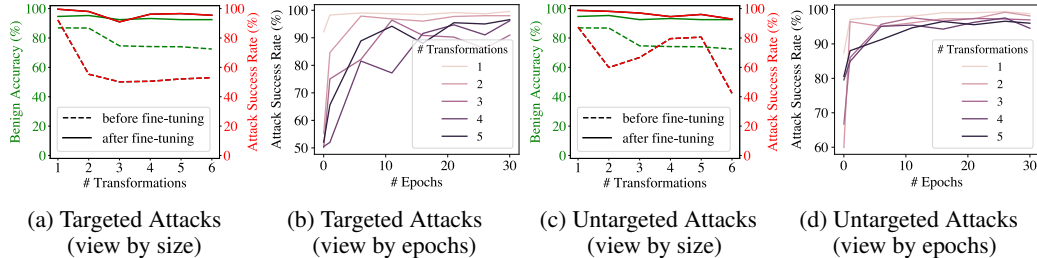


Figure 3: Performance of the BaRT defense on ImageNette with different numbers of transformations before and after fine-tuning the model. While the model achieves higher invariance, the defense becomes nearly ineffective², as evident from the top solid red curves in (a) and (c).

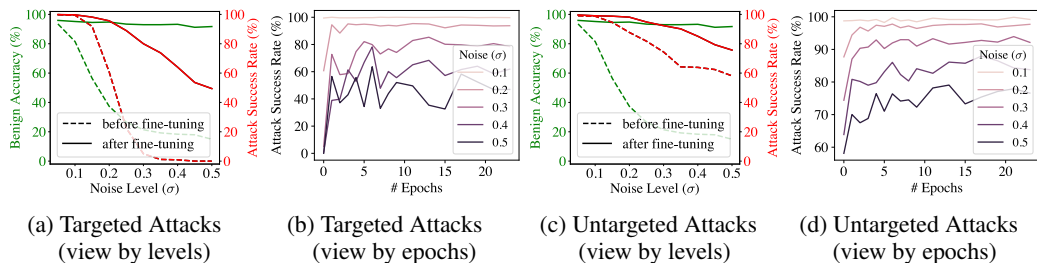


Figure 4: Performance of the randomized smoothing defense on ImageNette with different noise levels before and after fine-tuning the model. While the model achieves higher invariance, the defense becomes less effective³, as evident from the gap between dashed and solid red curves in (a) and (c).

270 invariance to their randomness. Here, we demonstrate this trade-off on realistic datasets, models, and
 271 defenses. In particular, we choose defenses with sufficient randomness (achieved in different ways)
 272 and compare their performance when being applied to models of different levels of invariance, where
 273 the invariance is achieved by applying the defense’s randomness to the training data so as to guide
 274 the model in learning their transformations.

275 **Randomness through Transformations.** We first examine the BaRT defense, which pre-processes
 276 input images with κ randomly composited stochastic transformations. It represents defenses aiming
 277 to increase randomness through diverse input transformations. Since our objective is to demonstrate
 278 the trade-off, it suffices to evaluate a subset of BaRT with $\kappa \leq 6$ transformations; this also avoids the
 279 training cost of evaluating the original BaRT with $\kappa = 25$. Figure 3 shows the performance of this
 280 defense with models before and after fine-tuning on its processed training data.

281 In Figures 3a and 3c, we first observe that fine-tuning indeed increases the model’s invariance to
 282 the applied defense’s randomness; the utility’s dashed green curves are improved to the solid green
 283 curves beyond 90%. However, as the model achieves higher invariance, the defense becomes nearly
 284 ineffective; the attack’s dashed red curves boost to the solid red curves near 100%. The same attack’s
 285 effectiveness throughout the fine-tuning procedure further verifies this observation, as shown in
 286 Figures 3b and 3d. It shows a clear trade-off between the defense’s robustness and the model’s
 287 invariance. That is, stochastic defenses start to lose robustness when their defended models achieve
 288 higher invariance to their transformations.

289 **Randomness through Noise Levels.** We then examine the randomized smoothing defense that adds
 290 Gaussian noise to the input image. Unlike BaRT’s diverse transformations, randomized smoothing
 291 increases randomness directly through the added noise’s variance σ^2 . This allows us to rigorously
 292 increase the randomness without unexpected artifacts like non-differentiable components. We evaluate
 293 the performance of this defense ($\sigma \leq 0.5$) with models before and after fine-tuning on training data
 294 perturbed with designated Gaussian noise. The results are shown in Figure 4.

²The defense may not grow stronger with more transformations, which is a drawback of BaRT that we will discuss in Appendix D.3. Yet, our evaluations focus on the fact that solid curves are above the dashed curves.

³One may also observe a trade-off between robustness and *utility* by examining the curve’s horizontal trend. However, we focus on the trade-off between robustness and *invariance*, which manifests in the vertical gap.

295 In Figures 4a and 4c, fine-tuning improves the model’s invariance, but the defense also becomes
296 significantly weaker during this process. For example, the targeted attack is nearly infeasible when
297 the model is variant to the large noise ($\sigma \geq 0.3$), yet is significantly more effective when the model
298 becomes invariant. The fine-tuning process in Figures 4b and 4d also verifies that stochastic defenses
299 become weaker when their defended models become more invariant to their randomness.

300 **Takeaways.** For both the BaRT and the randomized smoothing defense, we observe a clear trade-off
301 between the defense’s robustness and the model’s invariance to randomness. In particular, we find that
302 stochastic defenses lose adversarial robustness when their defended models achieve higher invariance
303 to their randomness. Our finding implies that such defenses would become ineffective when their
304 defended models are perfectly invariant to their randomness.

305 7 Discussions

306 In this section, we discuss several questions that arose from our study of the pre-processing defense.
307 The limitations and potential negative societal impacts of this work are discussed in Appendix F.

308 **What do pre-processing defenses really do?** We show that existing stochastic pre-processing
309 defenses do not introduce inherent robustness to the prediction task. Instead, they shift the input
310 distribution through randomness and transformations, which results in variance and introduces errors
311 during prediction. The observed “robustness”, in an unusual meaning for this literature, is a result of
312 these errors. This is fundamentally different from the inherent robustness provided by adversarial
313 training [22]. Although defenses like adversarial training still cost accuracy [33, 37], they do not
314 intentionally introduce errors like stochastic pre-processing defenses.

315 However, stochastic pre-processing defenses *do make* the attack harder when the adversary has only
316 limited knowledge of the defense’s transformations, e.g., in a low-query setting. In such a case, the
317 defense practically introduces noise to the attack’s optimization procedure, making it difficult for a
318 low-query adversary to find adversarial examples that consistently cross the probabilistic decision
319 boundary. In contrast, our theoretical analysis in Section 5 considers a powerful adversary with
320 full knowledge of the defense’s randomization space; hence it can optimize the adversarial example
321 directly towards the defended model’s decision boundary in expectation.

322 **How should we utilize randomness?** Stochastic defenses should rely on randomness that exploits
323 the properties of the prediction task. For example, some speech and vision problems are inherently
324 divisible into orthogonal subproblems (e.g., performing classification on portions of the speech
325 spectrum for keyword spotting [1]), where randomness can apply to the space of these subproblems.
326 Such defenses decouple robustness and invariance because (1) the models perform well on each
327 independent subproblem, and (2) attacks on one subproblem do not transfer to the other subproblems.
328 As such, an attacker has to target all possible subproblems, which reduces its effective attack budget.

329 **What are implications for adaptive attackers?** Our findings demonstrate that an adaptive attacker
330 needs to consider the spectrum of available standard attack algorithms, instead of just focusing on a
331 given attack algorithm because of the defense’s design. As we discover in this paper, EOT can be
332 unnecessary for seemingly immune stochastic defenses, yet its application to break these said defenses
333 gives a false impression about their security against weak attackers. When evaluating the robustness
334 of a defense, the adaptive attack should start by tuning standard approaches, before resorting to more
335 involved attack strategies. This approach helps us to identify the minimally capable attack that breaks
336 the defense and develop a better understanding of the defense’s fundamental weaknesses.

337 8 Conclusion

338 In this paper, we investigate stochastic pre-processing defenses and explain their limitations both
339 empirically and theoretically. We show that most stochastic pre-processing defenses are weaker than
340 previously thought, and recent defenses that indeed exhibit more randomness still face a trade-off
341 between their robustness and the model’s invariance to their transformations. While defending against
342 adversarial examples remains an open problem and designing proper adaptive evaluations is arguably
343 challenging, we demonstrate that stochastic pre-processing defenses are fundamentally flawed in their
344 current form. Our findings suggest that future work will need to find new ways of using randomness
345 that decouples the robustness and invariance.

346 **References**

- 347 [1] Shimaa Ahmed, Iliia Shumailov, Nicolas Papernot, and Kassem Fawaz. Towards more robust keyword
348 spotting for voice assistants. In *31st USENIX Security Symposium*, 2022.
- 349 [2] Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of
350 security: Circumventing defenses to adversarial examples. In Jennifer G. Dy and Andreas Krause, editors,
351 *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm*,
352 *Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages
353 274–283. PMLR, 2018. URL <http://proceedings.mlr.press/v80/athalye18a.html>.
- 354 [3] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples.
355 In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on*
356 *Machine Learning, ICML 2018, Stockholm*, *Stockholm, Sweden, July 10-15, 2018*, volume 80 of
357 *Proceedings of Machine Learning Research*, pages 284–293. PMLR, 2018. URL <http://proceedings.mlr.press/v80/athalye18b.html>.
- 359 [4] Battista Biggio, Iginio Corona, Davide Maiorca, Blaine Nelson, Nedim Srdic, Pavel Laskov, Giorgio Giac-
360 into, and Fabio Roli. Evasion attacks against machine learning at test time. In Hendrik Blockeel, Kristian
361 Kersting, Siegfried Nijssen, and Filip Zelezny, editors, *Machine Learning and Knowledge Discovery in*
362 *Databases - European Conference, ECML PKDD 2013, Prague, Czech Republic, September 23-27, 2013,*
363 *Proceedings, Part III*, volume 8190 of *Lecture Notes in Computer Science*, pages 387–402. Springer, 2013.
364 doi: 10.1007/978-3-642-40994-3_25. URL https://doi.org/10.1007/978-3-642-40994-3_25.
- 365 [5] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. In *2017*
366 *IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 39–57,
367 2017. doi: 10.1109/SP.2017.49. URL <https://doi.org/10.1109/SP.2017.49>.
- 368 [6] Jeremy M. Cohen, Elan Rosenfeld, and J. Zico Kolter. Certified adversarial robustness via random-
369 ized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th*
370 *International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California,*
371 *USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320. PMLR, 2019. URL
372 <http://proceedings.mlr.press/v97/cohen19c.html>.
- 373 [7] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of
374 diverse parameter-free attacks. In *Proceedings of the 37th International Conference on Machine Learning,*
375 *ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*,
376 pages 2206–2216. PMLR, 2020. URL <http://proceedings.mlr.press/v119/croce20b.html>.
- 377 [8] Guneet S. Dhillon, Kamyar Azizzadenesheli, Zachary C. Lipton, Jeremy Bernstein, Jean Kossaifi, Aran
378 Khanna, and Animashree Anandkumar. Stochastic activation pruning for robust adversarial defense. In *6th*
379 *International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 -*
380 *May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL [https://openreview.net/](https://openreview.net/forum?id=H1uR4GZRZ)
381 [forum?id=H1uR4GZRZ](https://openreview.net/forum?id=H1uR4GZRZ).
- 382 [9] FastAI. The imagenette dataset. URL <https://github.com/fastai/imagenette>.
- 383 [10] Chuan Guo, Mayank Rana, Moustapha Cissé, and Laurens van der Maaten. Countering adversarial images
384 using input transformations. In *6th International Conference on Learning Representations, ICLR 2018,*
385 *Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
386 URL <https://openreview.net/forum?id=SyJ7C1WCb>.
- 387 [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
388 In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA,*
389 *June 27-30, 2016*, pages 770–778, 2016. doi: 10.1109/CVPR.2016.90. URL [https://doi.org/10.](https://doi.org/10.1109/CVPR.2016.90)
390 [1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- 391 [12] Jörn-Henrik Jacobsen, Jens Behrmann, Richard S. Zemel, and Matthias Bethge. Excessive invariance
392 causes adversarial vulnerability. In *7th International Conference on Learning Representations, ICLR*
393 *2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL [https://openreview.net/](https://openreview.net/forum?id=BkfbpsAcF7)
394 [forum?id=BkfbpsAcF7](https://openreview.net/forum?id=BkfbpsAcF7).
- 395 [13] Sandesh Kamath, Amit Deshpande, Subrahmanyam Kambhampati Venkata, and Vineeth N Bala-
396 subramanian. Can we have it all? on the trade-off between spatial and adversarial robustness of
397 neural networks. In M. Ranzato, A. Beygelzimer, K. Nguyen, P. S. Liang, J. W. Vaughan, and
398 Y. Dauphin, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 27462–
399 27474. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper/2021/file/](https://proceedings.neurips.cc/paper/2021/file/e6ff107459d435e38b54ad4c06202c33-Paper.pdf)
400 [e6ff107459d435e38b54ad4c06202c33-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/e6ff107459d435e38b54ad4c06202c33-Paper.pdf).

- 401 [14] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International*
402 *Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference*
403 *Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- 404 [15] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 405 [16] Mathias Lécuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness
406 to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy, SP*
407 *2019, San Francisco, CA, USA, May 19-23, 2019*, pages 656–672. IEEE, 2019. doi: 10.1109/SP.2019.00044.
408 URL <https://doi.org/10.1109/SP.2019.00044>.
- 409 [17] Bai Li, Changyou Chen, Wenlin Wang, and Lawrence Carin. Certified adversarial robustness with additive
410 noise. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox,
411 and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference*
412 *on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,*
413 *Canada*, pages 9459–9469, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/335cd1b90bfa4ee70b39d08a4ae0cf2d-Abstract.html)
414 [335cd1b90bfa4ee70b39d08a4ae0cf2d-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/335cd1b90bfa4ee70b39d08a4ae0cf2d-Abstract.html).
- 415 [18] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *5th Interna-*
416 *tional Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference*
417 *Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=Skq89Scxx>.
- 418 [19] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference*
419 *on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
420 URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 421 [20] James Lucas, Shengyang Sun, Richard S. Zemel, and Roger B. Grosse. Aggregated momentum: Stability
422 through passive damping. In *7th International Conference on Learning Representations, ICLR 2019, New*
423 *Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Syxt5oC5YQ)
424 [id=Syxt5oC5YQ](https://openreview.net/forum?id=Syxt5oC5YQ).
- 425 [21] Clare Lyle, Mark van der Wilk, Marta Kwiatkowska, Yarin Gal, and Benjamin Bloem-Reddy. On the
426 benefits of invariance in neural networks, 2020.
- 427 [22] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards
428 deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Re-*
429 *presentations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*,
430 2018. URL <https://openreview.net/forum?id=rJzIBfZAb>.
- 431 [23] Tianyu Pang, Kun Xu, and Jun Zhu. Mixup inference: Better exploiting mixup to defend adversarial
432 attacks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia,*
433 *April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=ByxtC2VtPB>.
- 434 [24] Edward Raff, Jared Sylvester, Steven Forsyth, and Mark McLean. Barrage of random
435 transforms for adversarially robust defense. In *IEEE Conference on Computer Vision*
436 *and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages
437 6528–6537. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00669.
438 URL [http://openaccess.thecvf.com/content_CVPR_2019/html/Raff_Barrage_of_Random_](http://openaccess.thecvf.com/content_CVPR_2019/html/Raff_Barrage_of_Random_Transforms_for_Adversarially_Robust_Defense_CVPR_2019_paper.html)
439 [Transforms_for_Adversarially_Robust_Defense_CVPR_2019_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Raff_Barrage_of_Random_Transforms_for_Adversarially_Robust_Defense_CVPR_2019_paper.html).
- 440 [25] Kevin Roth, Yannic Kilcher, and Thomas Hofmann. The odds are odd: A statistical test for detecting
441 adversarial examples. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th*
442 *International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California,*
443 *USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 5498–5507. PMLR, 2019. URL
444 <http://proceedings.mlr.press/v97/roth19a.html>.
- 445 [26] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej
446 Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale
447 visual recognition challenge. *Int. J. Comput. Vis.*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
448 URL <https://doi.org/10.1007/s11263-015-0816-y>.
- 449 [27] Vasu Singla, Songwei Ge, Basri Ronen, and David Jacobs. Shift invariance can reduce adver-
450 sarial robustness. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman
451 Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 1858–
452 1871. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper/2021/file/](https://proceedings.neurips.cc/paper/2021/file/0e7c7d6c41c76b9ee6445ae01cc0181d-Paper.pdf)
453 [0e7c7d6c41c76b9ee6445ae01cc0181d-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/0e7c7d6c41c76b9ee6445ae01cc0181d-Paper.pdf).

- 454 [28] Chawin Sitawarin, Zachary Golan-Strieb, and David Wagner. Demystifying the adversarial robustness
455 of random transformation defenses. In *The AAAI-22 Workshop on Adversarial Machine Learning and*
456 *Beyond*, 2021. URL <https://openreview.net/forum?id=p4SrFydW05>.
- 457 [29] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and
458 Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun, editors, *2nd*
459 *International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014,*
460 *Conference Track Proceedings*, 2014. URL <http://arxiv.org/abs/1312.6199>.
- 461 [30] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking
462 the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern*
463 *Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2818–2826. IEEE Computer
464 Society, 2016. doi: 10.1109/CVPR.2016.308. URL <https://doi.org/10.1109/CVPR.2016.308>.
- 465 [31] Florian Tramèr, Jens Behrmann, Nicholas Carlini, Nicolas Papernot, and Jörn-Henrik Jacobsen. Fun-
466 damental tradeoffs between invariance and sensitivity to adversarial perturbations. In *Proceedings of*
467 *the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event,*
468 *volume 119 of Proceedings of Machine Learning Research*, pages 9561–9571. PMLR, 2020. URL
469 <http://proceedings.mlr.press/v119/tramer20a.html>.
- 470 [32] Florian Tramèr, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks
471 to adversarial example defenses. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-
472 Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems*
473 *33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, Decem-*
474 *ber 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/11f38f8ecd71867b42433548d1078e38-Abstract.html)
475 [11f38f8ecd71867b42433548d1078e38-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/11f38f8ecd71867b42433548d1078e38-Abstract.html).
- 476 [33] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Ro-
477 bustness may be at odds with accuracy. In *7th International Conference on Learning Representations,*
478 *ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=SyxAb30cY7>.
- 480 [34] Chang Xiao, Peilin Zhong, and Changxi Zheng. Enhancing adversarial defense by k-winners-take-all.
481 In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April*
482 *26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=Skgyv64tvr>.
- 483 [35] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan L. Yuille. Mitigating adversarial
484 effects through randomization. In *6th International Conference on Learning Representations, ICLR 2018,*
485 *Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
486 URL <https://openreview.net/forum?id=Sk9yuq10Z>.
- 487 [36] Yuzhe Yang, Guo Zhang, Zhi Xu, and Dina Katabi. Me-net: Towards effective adversarial robustness with
488 matrix estimation. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th*
489 *International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California,*
490 *USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 7025–7034. PMLR, 2019. URL
491 <http://proceedings.mlr.press/v97/yang19e.html>.
- 492 [37] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P. Xing, Laurent El Ghaoui, and Michael I. Jordan.
493 Theoretically principled trade-off between robustness and accuracy. In Kamalika Chaudhuri and Ruslan
494 Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019,*
495 *9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*,
496 pages 7472–7482. PMLR, 2019. URL <http://proceedings.mlr.press/v97/zhang19p.html>.

497 **Checklist**

- 498 1. For all authors...
- 499 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
500 contributions and scope? [Yes]
- 501 (b) Did you describe the limitations of your work? [Yes] See Appendix F.
- 502 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
503 Appendix F.
- 504 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
505 them? [Yes]
- 506 2. If you are including theoretical results...
- 507 (a) Did you state the full set of assumptions of all theoretical results? [Yes] We outline all
508 details for theoretical results in Appendix C.
- 509 (b) Did you include complete proofs of all theoretical results? [Yes] We include complete
510 proofs in Appendix C.
- 511 3. If you ran experiments...
- 512 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
513 imental results (either in the supplemental material or as a URL)? [Yes] Our code is
514 available in the supplementary material.
- 515 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
516 were chosen)? [Yes] We specify brief settings in Section 6.1 and provide complete
517 details in Appendix D.
- 518 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
519 ments multiple times)? [No] We study stochastic defenses, whose prediction procedure
520 is already a majority vote over multiple times.
- 521 (d) Did you include the total amount of compute and the type of resources used (e.g., type
522 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix D.
- 523 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 524 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 525 (b) Did you mention the license of the assets? [Yes]
- 526 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
527 We only included our own code for evaluation.
- 528 (d) Did you discuss whether and how consent was obtained from people whose data you’re
529 using/curating? [N/A] We only used public datasets as discussed in Section 6.1.
- 530 (e) Did you discuss whether the data you are using/curating contains personally identifiable
531 information or offensive content? [N/A] We did not use such data.
- 532 5. If you used crowdsourcing or conducted research with human subjects...
- 533 (a) Did you include the full text of instructions given to participants and screenshots, if
534 applicable? [N/A]
- 535 (b) Did you describe any potential participant risks, with links to Institutional Review
536 Board (IRB) approvals, if applicable? [N/A]
- 537 (c) Did you include the estimated hourly wage paid to participants and the total amount
538 spent on participant compensation? [N/A]