

# Increasing Confidence in Adversarial Robustness Evaluations

Roland Zimmermann\*  
University of Tübingen

Wieland Brendel  
University of Tübingen

Florian Tramèr  
Google

Nicholas Carlini  
Google

## Abstract

*Hundreds of defenses have been proposed in the past years to make deep neural networks robust against minimal (adversarial) input perturbations. However, only a handful of these could hold up their claims because correctly evaluating robustness is extremely challenging: Weak attacks often fail to find adversarial examples even if they unknowingly exist, thereby making a vulnerable network look robust. In this paper, we propose a test to identify weak attacks. Our test introduces a small and simple modification into a neural network that guarantees the existence of an adversarial example for every sample. Consequentially, any correct attack must succeed in attacking this modified network. For eleven out of thirteen previously-published defenses, the original evaluation of the defense fails our test, while stronger attacks that break these defenses pass it. We hope that attack unit tests such as ours will be a major component in future robustness evaluations and increase confidence in an empirical field that today is riddled with skepticism and disbelief. Online version & Code: [zimmerrol.github.io/active-tests/](https://zimmerrol.github.io/active-tests/)*

## 1. Introduction

Suppose that someone presents you with a purported proof that  $P \neq NP$ . The proof is long, complicated, and difficult to follow. How would you go about checking if this proof is correct?

One cumbersome way would be to directly refute the proof’s claim, e.g., to demonstrate that actually  $P=NP$  by designing an algorithm that solves 3-SAT in polynomial time. While this would definitely refute the proof, it is likely orders of magnitude more difficult than simply showing the proof is incorrect. Accordingly, researchers typically refute incorrect proofs by studying proofs line-by-line, until they identify some major flaw in the reasoning.

We argue a similar approach should be taken when evaluating adversarial example defenses. Evaluating defenses

to adversarial examples has proven to be extremely difficult [9]. In many areas of machine learning, evaluating the performance of a new technique is often trivial — for example by computing accuracy on some held-out test set. However evaluating defense robustness necessarily involves reasoning over *all* possible adversaries, and showing *none* can succeed. That is, a defense evaluation aims to prove that something is impossible. As a result, despite significant evaluation effort, most published defenses are quickly broken by stronger attacks [3, 9, 11, 14, 38].

This paper argues for viewing defense proposals as theorem statements, and the corresponding evaluations as proofs. The purpose of a defense evaluation, then, is to provide a convincing and rigorous argument that the defense is correct. Currently, for an adversary to claim to have a “break” of a defense, it is necessary to actually produce the adversarial examples that cause the model to make an error — analogous to refuting a complexity-theoretic impossibility result by producing an efficient algorithm. We argue that this is not how things should work. A valid refutation of a theorem would be to say “there is a flaw in your proof on line 9”. Because the null hypothesis for a theorem is that it is false, just as the null hypothesis for a defense should be that it is not robust.

Unfortunately, for defenses against adversarial examples, outside of studying the actual code used to implement the attack, there are relatively few opportunities to identify flawed evaluations by reading the paper. As a result, the current state-of-the-art in identifying flawed evaluations is to look for artifacts that indicate something has gone wrong — for example, that the attack fails even when it is allowed to construct *unbounded* perturbations) [9].

We develop a new *active robustness test* to complement existing (passive) tests [9, 26]. Our test designs a new task that is solvable by any sufficiently strong attack. Our test purposefully injects adversarial examples into a defense and then checks if the attack used to evaluate the defense is able to find them. If the attack fails this test, we know that it is too weak to distinguish between a robust and a non-robust defense, and thus the evaluation should not be trusted.

Our test would have potentially identified eleven out of

\*Work done while at Google. Contact: [roland.zimmermann@uni-tuebingen.de](mailto:roland.zimmermann@uni-tuebingen.de)

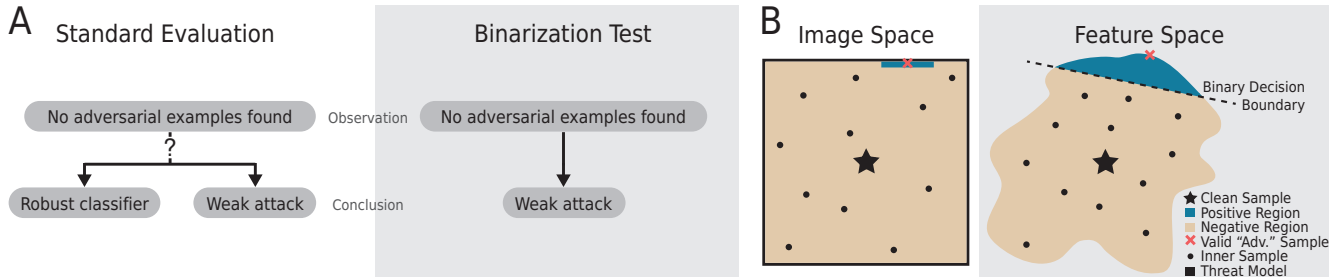


Figure 1. **A: Reasons for seemingly high robustness.** There are two reason an attack might not find an adversarial example. Either the classifier is robust or the attack is too weak and it could not find the existing adversarials. In our proposed *Binarization Test* we pose a new binary classification problem based on the original classifier such that adversarial examples always exist. Thus, if the attack does not find an adversarial example it follows that the attack is too weak. **B: Setup of the Binarization Test.** We construct a binary classification problem around a clean example such that there exists a valid “adversarial” example within the feasible set of the attack’s threat model. Based on the original classifier’s features, we create a new binary classifier whose robustness can be evaluated with the same evaluation function as used for the original classifier.

thirteen flawed evaluations found in previously-published papers. We hope that our testing methodology can become a standard component of future defense evaluations. To this end, defenses with exceptionally novel or different techniques, training algorithms, or architectures, may need to develop their own tailored version of our active unit test, in order to ensure the correctness of the defense evaluation.

## 2. Background

**Adversarial Examples** Adversarial examples contain imperceptible perturbations that change the decision of a deep neural network in arbitrary directions [4, 37]. Since they can manipulate the behavior of a model, they are seen as a security concern for machine learning applications. To find adversarial examples for a network one looks for inputs changing the output of the network while being close (under some norm) to the original data sample. There are a number of methods to solve this optimization problem and to attack a network. Adversarial attacks can be divided into white box methods that use gradient information about the model [e.g., 6, 11, 13, 14, 20], and black box methods that only use the output of the network [e.g., 1, 2, 5, 17, 23].

**Defenses** With an increasing awareness of the risk posed by adversarial examples, a vast number of defenses were proposed to increase adversarial robustness. For example, some defenses rely on additional input pre-processing [e.g., 16], some introduce architectural changes [e.g., 41], and others propose methods for detecting adversarial examples [e.g., 21]. However, most of these defenses eventually turned out to be ineffective against stronger attacks after publication [3, 38]. Until now only adversarial training [20] and its variants [e.g., 15, 27, 28] stood the test of time and could not be circumvented. A different approach to defend classifiers against adversarial perturbations are certified defenses which give a theoretical guarantee of the classifier’s

robustness. Yet, the robustness of these approaches does not yet reach that of adversarial training [12, 19, 40].

**Challenges in Evaluating Defenses** Properly evaluating the robustness of a model against adversarial examples is non-trivial and there are many potential pitfalls [9]. The critical issue is that when a defense is shown to be robust to a specific attack, this either means that the model is truly robust, or that the attack is suboptimal (see Figure 1A). Possible reasons for an attack to be ineffective are incorrect hyperparameters, or mechanisms in the model that (unintentionally) hinder the attack’s optimization process [3]. Examples include defenses built around non-continuous activation functions [e.g., 41] or relying on vanishing gradients [e.g., 36]. To address the former issue, prior work has developed attacks that alleviate the need to manually tune hyperparameters [14]. But as we will show, these attacks are not guaranteed to work well for any model. While the latter issue can be counteracted by using adaptive attacks [38] that are adjusted to a specific model’s idiosyncrasies, it remains non-trivial to detect suboptimal attacks in the first place. Previous work suggested guidelines for evaluating the adversarial robustness of a model [9] or developed (passive) indicator values hinting at a failed evaluation [26]. These indicator values are based on metrics tracked during an adversarial attack and check for certain failure cases. Our work goes beyond these indicator values by arguing that researchers should *actively* demonstrate their adversarial attack works and is sufficiently strong, and that their empirical findings can be trusted.

## 3. Active Attack Evaluation Tests

The evaluation of a defense against adversarial attacks becomes more reliable — and the estimated robustness more correct — if the attack is believed to be sufficiently strong. The strength of an attack is not an absolute value

but depends on the defense it is meant to evaluate, as various defense mechanisms hinder specific attacks [3]. Thus, for a new defense one needs to demonstrate that the attack proposed to evaluate it is appropriate. In this section, we propose a test that measures the adequacy of a defense’s evaluation scheme, and is thereby able to warn researchers of potentially unreliable robustness claims.

As stated before, to empirically demonstrate the robustness of a classifier  $f$  for some input  $\mathbf{x}_c$  one runs an attack and shows that it fails to find an adversarial example  $\mathbf{x}_{adv}$  within distance  $d(\mathbf{x}_c, \mathbf{x}_{adv}) \leq \epsilon$ . But can one really be sure there are no adversarial examples in the  $\epsilon$  ball if the attack fails? Since the attack cannot give a guarantee for this, there might still be stronger attacks that do find adversarial examples (see Figure 1A).

We propose a test that enables researchers to check whether an attack is too weak to support their robustness claims. In our test we craft a new classifier that is as similar as possible to the original, but where we intentionally inject an adversarial example  $\mathbf{x}_{adv}$ . Then, we measure the robustness of the new (by definition vulnerable) classifier by running the evaluation method and checking whether an adversarial example is found. If the originally used attack fails to find adversarial examples for the modified classifier, we cannot expect it to properly estimate the robustness of the original classifier either.

### Test for Classifiers with Linear Classification Readouts

We begin by describing how to construct the modified vulnerable model for a classifier  $f$  that consists of a feature extractor  $f^*$  followed by a linear classification head. Any standard neural network architecture falls into this category: the feature extractor  $f^*$  is every layer except the last, and the linear classification head is the final logit projection layer. We keep the feature extractor  $f^*$  unchanged to avoid changing the fundamental behavior of the model, but replace the classification readout with a newly trained module. This module is trained on a new, specially constructed dataset. This dataset allows us to reliably create a classifier where—by design—there exists at least one adversarial example for each sample. A pseudocode definition of our test is shown in Algorithm 1. In detail, for each test sample  $\mathbf{x}_c$  our test consists of the following steps:

Initially, we create two collections of input samples which are perturbed versions of  $\mathbf{x}_c$

$$\mathcal{X}_i := \{ \hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) < \xi \cdot \epsilon \wedge \hat{\mathbf{x}} \neq \mathbf{x}_c \} \cup \{ \mathbf{x}_c \}_{1, \dots, N_i} \text{ and}$$

$$\mathcal{X}_b := \{ \hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) = \epsilon \}_{1, \dots, N_b},$$

which are sets of points from the inside and the boundary of the  $\epsilon$ -ball, respectively, with size  $N_i, N_b > 0$ . Further,  $\xi \in (0, 1)$  controls the margin between the inner  $\mathcal{X}_i$  and the boundary set  $\mathcal{X}_b$ . Decreasing  $\xi$  effectively increases the gap

---

### Algorithm 1 Binarization Test for classifiers with linear classification readouts

---

**input:** test samples  $\mathcal{X}_{test}$ , feature extractor  $f^*$  of original classifier, number of inner/boundary samples  $N_i$  and  $N_b$ , distance  $\epsilon$ , sampling functions for data from the inside/boundary of the  $\epsilon$ -ball.

```

function BINARIZATIONTEST( $f^*$ ,  $\mathcal{X}_{test}$ ,  $N_b$ ,  $N_i$ ,  $\epsilon$ )
  attack_success = []
  rnd_attack_success = []
  for all  $\mathbf{x}_c \in \mathcal{X}_{test}$  do
     $b = \text{CreateBinaryClassifier}(f^*, \mathbf{x}_c, \epsilon)$ 
    # evaluate robustness of binary classifier
    attack_success.insert( $\text{RunAttack}(b, \mathbf{x}_c)$ )
    rnd_attack_success.insert( $\text{RunRndAttack}(b, \mathbf{x}_c)$ )
  ASR = Mean(attack_successful)
  RASR = Mean(random_attack_successful)
  return ASR, RASR
end function

```

```

function CREATEBINARYCLASSIFIER( $f^*$ ,  $\mathbf{x}_c$ )
  # draw input samples around clean example
   $\mathcal{X}_i = \{ \mathbf{x}_c \} \cup \{ \text{SampleInnerPoint}(\mathbf{x}_c, \epsilon) \}_{1, \dots, N_i}$ 
   $\mathcal{X}_b = \{ \text{SampleBoundaryPoint}(\mathbf{x}_c, \epsilon) \}_{1, \dots, N_b}$ 
  # get features for images
   $\mathcal{F}_i = \{ f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_i \}$ 
   $\mathcal{F}_b = \{ f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_b \}$ 
  # define labels & create labeled dataset
   $\mathcal{D} = \{ (\hat{\mathbf{x}}, 0) \mid \hat{\mathbf{x}} \in \mathcal{F}_i \} \cup \{ (\hat{\mathbf{x}}, 1) \mid \hat{\mathbf{x}} \in \mathcal{F}_b \}$ 
  # train linear readout on extracted features
   $b = \text{TrainReadout}(\mathcal{D})$ 
  return binary classifier  $b$  based on feature encoder  $f^*$ 
end function

```

---

between inner and boundary points, thus, making it easier to distinguish between the two sets of samples.

Next, for every sample in each of the two sets, obtain the feature representation of the penultimate layer of  $f$ ,

$$\mathcal{F}_i := \{ f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_i \} \text{ and } \mathcal{F}_b := \{ f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_b \}$$

Now, train a linear (binary) discriminator  $g$  that distinguishes samples from  $\mathcal{F}_i$  and  $\mathcal{F}_b$ , i.e., it distinguishes between mildly perturbed images — the interior of the  $\epsilon$ -ball — and some more strongly perturbed images — on the boundary of the  $\epsilon$ -ball. To mimic the behavior of the normal classifier’s readout as much as possible, we roughly match the value range of the predicted logits. We want to make sure there exists at least one sample within the threat model’s  $\epsilon$ -ball that  $g$  classifies differently than the original sample. Thus, we need to ensure that  $g$  achieves a perfect accuracy on these two sets. If this is not possible for a sample  $\mathbf{x}_c$ , we cannot apply the test and, hence, skip the sample.

By combining the original classifier’s feature extractor  $f^*$  with the binary discriminator  $g$ , i.e.  $h = g \circ f^*$ , one gets a new classifier that maps samples to a binary decision. Most importantly, each boundary sample  $\mathcal{X}_b$  acts as an  $\epsilon$ -bounded adversarial example  $\mathbf{x}_{adv}$  for the clean sample  $\mathbf{x}_c$ .

We are interested in two properties of this classifier  $h$ :

1. The efficacy of the used evaluation method/adversarial attack. For this, one uses the original adversarial attack to attack the modified model  $h$  for the clean sample  $\mathbf{x}_c$  and records whether an adversarial sample  $\mathbf{x}^*$  within the allowed  $\epsilon$ -ball is found. When calculated and averaged over multiple samples, we call this value the *test score*.
2. The difficulty of the test. To assess this, we use a model-agnostic attack, namely a purely randomized one. We attack the modified classifier  $h$  by randomly sampling approximately as many additional data points from within the  $\epsilon$ -ball around the clean sample  $\mathbf{x}_c$  as the adversarial attack queries the model, e.g. for an  $N$ -step PGD attack [20] use  $N$  additional random samples. Finally, one tests whether at least one of them turns out to be an adversarial perturbation for  $h$ . By averaging over multiple samples, we get the *random attack success rate* (R-ASR).

Note that if the classifier  $f$  does not use a linear classification readout, one has to modify the test slightly: Instead of using a linear readout for  $g$  one needs to use the same type of mechanism that was used originally. While this modification is conceivable for various mechanisms, e.g.  $k$ -nearest neighbors classification [35], there might be architectures for which this is not possible, e.g., classification through likelihood estimations based on generative models [32, 45].

**Test for Models Leveraging Detectors** Appendix A shows how to adapt this test to detection defenses.

## 4. Evaluation

Our test would have potentially prevented the publication of eleven out of thirteen broken defenses. We apply our binarization test to thirteen defenses and show that it would have identified flaws in nine previously peer-reviewed (and then later broken) and two published (but not yet broken) defenses. All of these defenses assume an  $\ell_\infty$  threat model. The specific design choices for the tests adapted to each defense can be found in Appendix B.

**Defenses without Detectors** We analyze eight defenses which use a classifier with a linear classification readout [8, 22, 25, 31, 33, 41, 43, 44].

We additionally apply our test to two defenses that do not use a simple linear readout to perform classification. It is straightforward to adapt the binarization test defined in Algorithm 1 for these classifier architectures. The defense

by Verma et al. [39] leverages an ensemble of readouts. For our test we therefore also train an ensemble of binary readouts. The classifier of Pang et al. [24] learns to map images to pre-defined class-prototype vectors, and then uses nearest neighbour classification. We reflect this in the test by using two of the class prototypes an associate them with the inner and boundary samples, respectively. Then we re-train a linear layer mapping from features to class prototypes.

In fact, we are the first to show that the defense by Sarkar et al. [31] is less robust than originally reported, as suggested by the fact that it fails our test. All of the above defenses except those by Sarkar et al. [31] were known to be flawed and have been circumvented before.

**Defenses with Detectors** We investigate three published defenses that aim to detect adversarial perturbations. Following previous work [7], we analyze each defense in a setting where it achieves a false positive rate of 5%. While the detection algorithm proposed by Roth et al. [29] runs statistical tests on the classifier’s confidence, Shan et al. [34] and Yang et al. [42] analyze earlier activations of the classifier. The first two defenses have been broken before [7, 38] while the latter had not been independently re-evaluated.

### 4.1. Evaluation of Not-previously-broken Defenses

We begin by investigating the two recent and not yet broken defenses. Here, we are interested in seeing whether the original robustness evaluations pass our binarization test. While a positive result would increase confidence in the defenses’ claims, a negative outcome would cast doubts.

**Sarkar et al. [31]** The original evaluation of this defense fails our test with a test score of 0.04. This is strong evidence that the attack is weak and thus the robustness claim likely overestimated. Upon investigation, we found a flaw in the original evaluation’s code: The statistics of the batch normalization layers are not frozen during evaluation, which changes the behavior of the model during the attack. Properly freezing these layers at inference and increasing the number of PGD steps from 20 to 75 yields a perfect score (1.0) in our binarization test. Moreover, this updated evaluation methodology reduces the robust accuracy to  $\leq 1\%$  down from the originally reported 60.15% and, thus, effectively breaks the defense.

**Yang et al. [42]** For this detector-based defense, we find that the attack used in the original evaluation is agnostic to the detector and only targets the classifier. Consequentially, this attack fails our binarization test with a low score of 0.26. We thus create a new adversarial attack based on PGD that combines two objectives: (1) fool the classifier by maximizing the adversarial loss and (2) stay undetected by matching the features of a non-adversarial sample as much

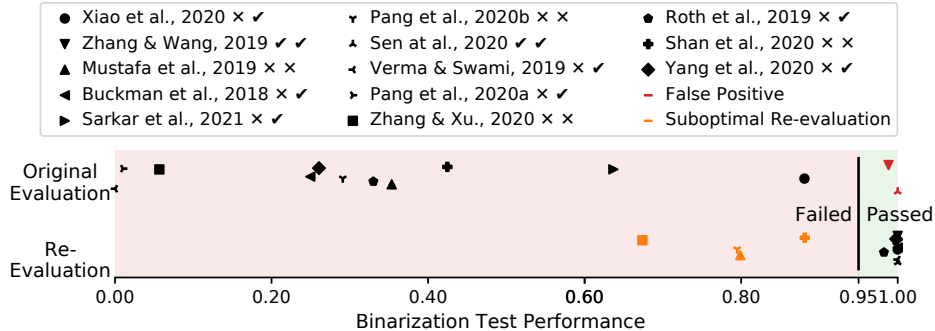


Figure 2. **The binarization test identifies flawed adversarial evaluations.** The x-axis shows the score in our proposed binarization test for the original attack (upper) and a subsequent improved attack (lower). We define a threshold of 0.95 that attacks need to achieve to pass our test. Note that for each defense, the improved attack substantially decreases the defense’s robust accuracy (by at least 12%). Black markers indicate original attacks that fail the test, as well as improved attacks that pass the test (i.e., true positives and true negatives for our test). Red markers indicate suboptimal original evaluations that still pass our test (false positives). Orange markers indicate re-evaluations that used suboptimal attacks (as shown by our test) that still broke the defense. We discuss these cases in Section 4.2. Checks and crosses in the legend indicate passing/failing tests for the original and the re-evaluation, respectively. See Appendix, Figure 4 for the robust accuracies.

as possible (a feature matching attack [30, 38]). This adaptive attack achieves a nearly perfect score of 0.99 in the test and reduces the robust accuracy of the defense from the originally reported 99 % down to  $\leq 12$  %.

#### 4.2. Interpreting Test Results for Weak & Strong Attacks

Since eleven of the considered defenses have already been broken before, and we showed how to break the remaining two, we now have access to both a flawed and a well-working adversarial evaluation method for each defense. This allows us to compare how these attacks perform in terms of both the estimated robust accuracy and the score on the binarization test. We visualize the results in Figure 2. For eleven out of the thirteen considered defenses, our proposed test would have flagged their evaluation as insufficient: the original attacks’ test performance is substantially below a perfect score. Furthermore, the test scores improve for almost all defenses when replacing the originally used evaluation code with an improved attack.

**Explaining the False Positives** Our test incorrectly lets two defense evaluations that had bugs pass (see red markers in Figure 2). When investigating these failure cases in more detail we find that the original attack used by Sen et al. [33] is not bad or incorrectly implemented per se, but is not used correctly. Namely, the attack generates adversarial examples with respect to the classifier’s predicted label, instead of the ground-truth label. As a result, for some misclassified samples the attack actually *corrects* the classifier’s mistake! By switching to an attack that correctly targets the ground-truth label, we reduce the robust accuracy drastically.

Unfortunately, our test is not suited for catching such a mistake. Indeed, by design, our test constructs a binary

classifier with 100% accuracy (and thus the classifier’s predicted label is always equal to the ground-truth). If we view our proposed test as a *unit test* for an attack, then the type of bug in the above evaluation is akin to an *integration bug*, where the (correct) attack is incorrectly called.

For the defense by Zhang et al. [43] we notice a high R-ASR value ( $> 0.75$ ) that we could not decrease further. We hypothesize that by increasing the number of inner samples  $N_i$  substantially, the test might become hard enough to indicate sub-optimal evaluations for this defense.

**Explaining the Suboptimal Re-evaluations** There are also four defenses for which the improved attacks still fail our test, even though their test performance is better than for the original attacks. The authors of the improved attack for Pang et al. [25] note that while this attack already breaks the defense, one could improve the attack further. For the defenses by Mustafa et al. [22] and Zhang et al. [44] the improved attacks are not adaptive attacks but part of AutoAttack’s attack collection [14]. Although these attacks were sufficient to drastically reduce the measured robustness of the defenses (see Appendix, Figure 4), they are not guaranteed to be the optimal attacks for these defenses. While the attack [7] used for the re-evaluation breaks the defense by Shan et al. [34], the imperfect test score hints at an even more potent and yet-to-be-discovered adaptive attack.

#### 4.3. Hardness of the Test

To put the performance that an adversarial attacks achieves in the binarization test into perspective, we quantify the hardness of the test using the prior random attack success rate (R-ASR). Comparing it to the test result of the attack allows us to deduce how effective the attack is in finding adversarial examples for the model in question.

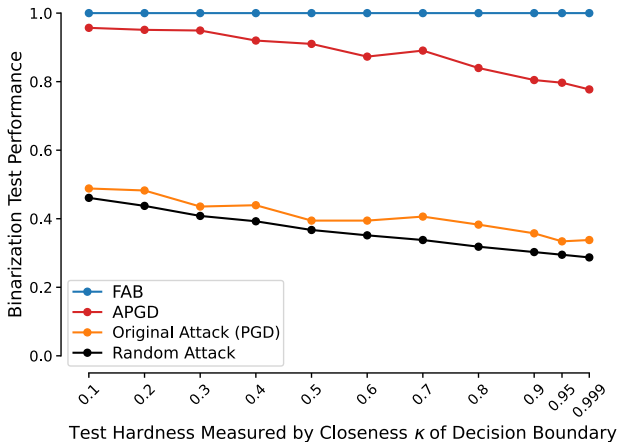


Figure 3. **Hyperparameters influence test’s hardness.** For the defense by Mustafa et al. [22], we compare the test performance of two sub-optimal attacks, namely the original PGD attack (orange) and AutoPGD (red) [14], yielding robust accuracies of 32.32 % and 8.16 % with that of the more optimal FAB attack (blue) [13] yielding 0.71 %. As one indicator of the test’s hardness, we show the ASR of a random attacker (R-ASR, black). Also, the test’s hardness is quantified by  $\kappa$  which, in feature space, measures the distance between decision boundary and boundary sample relative to the distance between boundary and closest inner sample.

There are several parameters and design choices relevant for our test that influence its hardness. For one, by increasing  $N_i$  we train the binary discriminator on a larger number of different non-adversarial points which increases robustness of the discriminator and, thus, makes the test harder. Contrary, by increasing  $N_b$  we plant a larger number of adversarial examples for the discriminator within the  $\epsilon$ -ball, making the test simpler.

Even with a large but finite number of training samples, there is no unique solution for the binary discriminator but instead a set of valid solutions. While all of these classifiers have perfect accuracy on the training set, they differ in how close the decision boundary is placed to the boundary samples. The closer the decision boundary to boundary samples, the smaller the volume of valid adversarial examples and, thus, the harder the test becomes. The effect of the decision boundary’s closeness on the test’s hardness is visualized in Figure 3 for example defense. Here, placing the boundary closer to the boundary samples decreases both the R-ASR as well as the ASR of two sub-optimal attacks while that of a better suited attack stays robustly at 1.0.

On the one hand a test that is too easy has no predictive power about the attack’s efficacy (since any attack might trivially pass it), while on the other hand a test that is too challenging might actually underestimate the attack’s true performance. Therefore, one needs to tune the test’s hardness to a reasonable level. For tweaking the test’s hard-

ness we recommend the following procedure: To make sure one does not overestimate (since this is the more dangerous direction) the test performance, start with a configuration that makes the test as hard as possible while still being computationally feasible. Now, decrease the hardness until the adversarial attack in question reaches an (almost) perfect ASR. Note that if there is no configuration that yields this, then the attack did not pass the test and one should be skeptical of the attack’s ability to properly estimate the classifier’s robustness. Finally, compare the ASR with the R-ASR: If the ASR is not substantially higher — or is even lower — than the R-ASR, this is strong evidence that the attack performs poorly. If instead the gap is large the attack has passed this necessary test and might be powerful enough to properly estimate the classifier’s robustness.

## 5. Discussion & Conclusion

This paper made a case for the need for *active* tests. The goal of an active test is to provide compelling evidence that an attack has sufficient power to evaluate a classifier’s robustness. We presented such a test for defenses using linear classification readouts and showed how to adapt this test for different defense mechanisms such as detector-based defenses. The type of test proposed in this work acts as a necessary condition for robustness evaluations, i.e., an attack that fails the test will most likely overestimate the classifier’s robustness.

While we have presented a potential test that could help defense authors demonstrate sufficient power of their adversarial evaluation, our tests cannot be comprehensive and apply to every possible defense. For example, all of our tests are primarily designed to work for defenses that use linear classification readout layers. If a defense were to have a different classification layer instead, such as a k-Nearest Neighbor classifier, then the tests we develop would not apply directly and need to be modified accordingly. Consequently, defense authors should aim to develop their own tests, depending on the particular claims that are made.

As we showed that this type of test would have prevented the publication of thirteen flawed defenses, we are optimistic that active tests can improve the reliability of future publications in the field of adversarial robustness.

## Acknowledgements

We thank Alexey Kurakin for his valuable feedback. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting RSZ. This work was supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039A. WB acknowledges financial support via an Emmy Noether Grant funded by the German Research Foundation (DFG) under grant no. BR 6382/1-1. WB is a member of the Machine Learning Cluster of Excellence, EXC number 2064/1 – Project number 390727645.

## References

- [1] Moustafa Alzantot, Yash Sharma, Supriyo Chakraborty, Huan Zhang, Cho-Jui Hsieh, and Mani B Srivastava. Genattack: Practical black-box attacks with gradient-free optimization. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pages 1111–1119, 2019. 2
- [2] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square attack: a query-efficient black-box adversarial attack via random search. In *European Conference on Computer Vision*, pages 484–501. Springer, 2020. 2
- [3] Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 274–283. PMLR, 2018. 1, 2, 3
- [4] Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrđić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *Joint European conference on machine learning and knowledge discovery in databases*, pages 387–402. Springer, 2013. 2
- [5] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. 2
- [6] Wieland Brendel, Jonas Rauber, Matthias Kümmeler, Ivan Ustyuzhaninov, and Matthias Bethge. Accurate, reliable and fast robustness evaluation. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 12841–12851, 2019. 2
- [7] Oliver Bryniarski, Nabeel Hingun, Pedro Pachuca, Vincent Wang, and Nicholas Carlini. Evading adversarial example detection defenses with orthogonal projected gradient descent. *arXiv preprint arXiv:2106.15023*, 2021. 4, 5
- [8] Jacob Buckman, Aurko Roy, Colin Raffel, and Ian J. Goodfellow. Thermometer encoding: One hot way to resist adversarial examples. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. 4
- [9] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness. *ArXiv preprint*, abs/1902.06705, 2019. 1, 2
- [10] Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pages 3–14, 2017. 11
- [11] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017. 1, 2
- [12] Jeremy M. Cohen, Elan Rosenfeld, and J. Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320. PMLR, 2019. 2
- [13] Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 2196–2205. PMLR, 2020. 2, 6
- [14] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 2206–2216. PMLR, 2020. 1, 2, 5, 6
- [15] Sven Gowal, Sylvestre-Alvise Rebuffi, Olivia Wiles, Florian Stimberg, Dan Andrei Calian, and Timothy A Mann. Improving robustness using generated data. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [16] Chuan Guo, Mayank Rana, Moustapha Cissé, and Laurens van der Maaten. Countering adversarial images using input transformations. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. 2
- [17] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 2142–2151. PMLR, 2018. 2
- [18] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 11
- [19] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. Certified robustness to adversarial examples with differential privacy. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 656–672. IEEE, 2019. 2
- [20] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. 2, 4
- [21] Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference*

- Track Proceedings*. OpenReview.net, 2017. 2, 10
- [22] Aamir Mustafa, Salman H. Khan, Munawar Hayat, Roland Goecke, Jianbing Shen, and Ling Shao. Adversarial defense by restricting the hidden space of deep neural networks. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 3384–3393. IEEE, 2019. 4, 5, 6
- [23] Nina Narodytska and Shiva Prasad Kasiviswanathan. Simple black-box adversarial attacks on deep neural networks. In *CVPR Workshops*, volume 2, 2017. 2
- [24] Tianyu Pang, Kun Xu, Yinpeng Dong, Chao Du, Ning Chen, and Jun Zhu. Rethinking softmax cross-entropy loss for adversarial robustness. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 4
- [25] Tianyu Pang, Kun Xu, and Jun Zhu. Mixup inference: Better exploiting mixup to defend adversarial attacks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 4, 5
- [26] Maura Pintor, Luca Demetrio, Angelo Sotgiu, Giovanni Manca, Ambra Demontis, Nicholas Carlini, Battista Biggio, and Fabio Roli. Indicators of attack failure: Debugging and improving optimization of adversarial examples. *ArXiv preprint*, abs/2106.09947, 2021. 1, 2
- [27] Rahul Rade and Seyed-Mohsen Moosavi-Dezfooli. Helper-based adversarial training: Reducing excessive margin to achieve a better accuracy vs. robustness trade-off. In *ICML 2021 Workshop on Adversarial Machine Learning*, 2021. 2
- [28] Sylvestre-Alvise Rebuffi, Sven Gowal, Dan Andrei Calian, Florian Stimberg, Olivia Wiles, and Timothy A Mann. Data augmentation can improve robustness. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [29] Kevin Roth, Yannic Kilcher, and Thomas Hofmann. The odds are odd: A statistical test for detecting adversarial examples. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 5498–5507. PMLR, 2019. 4
- [30] Sara Sabour, Yanshuai Cao, Fartash Faghri, and David J. Fleet. Adversarial manipulation of deep representations. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. 5, 10
- [31] Anindya Sarkar, Anirban Sarkar, Sowrya Gali, and Vineeth N Balasubramanian. Get fooled for the right reason: Improving adversarial robustness through a teacher-guided curriculum learning approach. *ArXiv preprint*, abs/2111.00295, 2021. 4
- [32] Lukas Schott, Jonas Rauber, Matthias Bethge, and Wieland Brendel. Towards the first adversarially robust neural network model on mnist. In *International Conference on Learning Representations*, 2018. 4
- [33] Sanchari Sen, Balaraman Ravindran, and Anand Raghunathan. EMPIR: ensembles of mixed precision deep networks for increased robustness against adversarial attacks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 4, 5
- [34] Shawn Shan, Emily Wenger, Bolun Wang, Bo Li, Haitao Zheng, and Ben Y Zhao. Gotta catch'em all: Using honeypots to catch adversarial attacks on neural networks. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, pages 67–83, 2020. 4, 5, 11
- [35] Chawin Sitawarin and David Wagner. Defending against adversarial examples with k-nearest neighbor. *arXiv preprint arXiv:1906.09525*, 2019. 4
- [36] Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. 2
- [37] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In Yoshua Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. 2
- [38] Florian Tramer, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. *Advances in Neural Information Processing Systems*, 33, 2020. 1, 2, 4, 5
- [39] Gunjan Verma and Ananthram Swami. Error correcting output codes improve probability estimation and adversarial robustness of deep neural networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 8643–8653, 2019. 4
- [40] Eric Wong and J. Zico Kolter. Provable defenses against adversarial examples via the convex outer adversarial polytope. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 5283–5292. PMLR, 2018. 2
- [41] Chang Xiao, Peilin Zhong, and Changxi Zheng. Enhancing adversarial defense by k-winners-take-all. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 2, 4
- [42] Puyudi Yang, Jianbo Chen, Cho-Jui Hsieh, Jane-Ling Wang, and Michael Jordan. MI-loo: Detecting adversarial examples with feature attribution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6639–6647, 2020. 4
- [43] Haichao Zhang and Jianyu Wang. Defense against adversarial attacks using feature scattering-based adversarial training. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett,



editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 1829–1839, 2019. [4](#), [5](#), [11](#)

- [44] Haichao Zhang and Wei Xu. Adversarial interpolation training: A simple approach for improving model robustness, 2020. [4](#), [5](#), [11](#)
- [45] Roland S Zimmermann, Lukas Schott, Yang Song, Benjamin A Dunn, and David A Klindt. Score-based generative classifiers. *arXiv preprint arXiv:2110.00473*, 2021. [4](#)

## A. Tests for Models Leveraging Detectors

A test similar to the one presented in Section 3 can also be used to validate the evaluation of detection defenses. These use an additional algorithm that detects and rejects adversarial examples [21]. As earlier, we assume that the classifier can be divided into a feature encoder and linear readout.

We define two tests (a “normal” and “inverted” test). Any reliable evaluation method must pass both. A pseudocode definition of the proposed tests is given as Algorithm 2.

**Normal Test** Adversarial examples for a detection defense need to change the classifier’s output while remaining undetected. Thus, we need to change the construction of the binary classifier slightly. Namely, we modify the set  $\mathcal{X}_b$  such that none of these samples gets rejected by the detector - in practice, we enforce this using rejection sampling, by redrawing boundary points until we find one that is undetected. Note, that we make no modifications to the detector, since this might require non-trivial optimization of the detector’s parameters. Some adversarial attacks (e.g., feature matching [30]) for detector defenses assume access to reference data samples that belong to a different class but are not

---

**Algorithm 2** Binarization Test for classifiers with a linear classification readout and a detector

---

**input:** test samples  $\mathcal{X}_{\text{test}}$ , feature extractor  $f^*$  of original classifier, adversarial detector  $d(\cdot)$  returning 1 for detected samples and 0 otherwise, number of inner/boundary/reference samples  $N_i/N_b/N_r$ , distance  $\epsilon$ , sampling functions for data from the inside/boundary of the  $\epsilon$ -ball, relative distance (in terms of  $\epsilon$ ) of positive and reference samples  $\eta > 1$ .

**function** BINARIZATIONTEST( $f^*, d, \mathcal{X}_{\text{test}}, N_b, N_i, N_r, \epsilon, \eta$ )

  attack\_success = []

  rnd\_attack\_success = []

**for all**  $\mathbf{x}_c \in \mathcal{X}_{\text{test}}$  **do**

$b, \mathcal{X}_r = \text{CreateBinaryClassifier}(f^*, \mathbf{x}_c, \epsilon)$

    # evaluate robustness of binary classifier

    attack\_success.insert( $\text{RunAttack}(b, d, \mathbf{x}_c, \mathcal{X}_r)$ )

    rnd\_attack\_success.insert( $\text{RunRndAttack}(b, d, \mathbf{x}_c)$ )

  ASR = Mean(attack\_successful)

  RASR = Mean(random\_attack\_successful)

**return** ASR, RASR

**end function**

**function** INVERTEDBINARIZATIONTEST( $f^*, d, \mathcal{X}_{\text{test}}, N_b, N_i, N_r, \epsilon, \eta$ )

  #  $\neg d$  denotes the negated/inverted detector

**return** BinarizationTest( $f^*, \neg d, \mathcal{X}_{\text{test}}, N_b, N_i, N_r, \epsilon, \eta$ )

**end function**

**function** CREATEBINARYCLASSIFIER( $f^*, \mathbf{x}_c, d$ )

  # draw input samples around clean example

$\mathcal{X}_i = \{ \mathbf{x}_c \} \cup \{ \text{SampleInnerPoint}(\mathbf{x}_c, \epsilon) \}_{1, \dots, N_i}$

$\mathcal{X}_b = \{ \text{SampleBoundaryPoint}(\mathbf{x}_c, \epsilon), d(z) = 1 \}_{1, \dots, N_b}$

  # get positive samples outside the  $\epsilon$ -ball, e.g., as a reference for logit matching attacks

$\mathcal{X}_r = \{ \text{SampleBoundaryPoint}(\mathbf{x}_c, \eta\epsilon), d(z) = 1 \}_{1, \dots, N_r}$

  # get features for images

$\mathcal{F}_i = \{ f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_i \}$

$\mathcal{F}_b = \{ f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_b \}$

$\mathcal{F}_r = \{ f^*(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}_r \}$

  # define labels & create labeled dataset

$\mathcal{D} = \{ (\hat{\mathbf{x}}, 0) \mid \hat{\mathbf{x}} \in \mathcal{F}_i \} \cup \{ (\hat{\mathbf{x}}, 1) \mid \hat{\mathbf{x}} \in \mathcal{F}_b \} \cup \{ (\hat{\mathbf{x}}, 1) \mid \hat{\mathbf{x}} \in \mathcal{F}_r \}$

  # train linear readout on extracted features

$b = \text{TrainReadout}(\mathcal{D})$

**return** binary classifier  $b$  based on feature encoder  $f^*$  and reference samples  $\mathcal{X}_r$

**end function**

---

adversarial and, thus, do not get rejected. In our setting this can be realized by randomly sampling data points outside the  $\epsilon$  ball. Thus, we create a new collection

$$\mathcal{X}_r := \{ \hat{\mathbf{x}} \mid d(\mathbf{x}_c, \hat{\mathbf{x}}) = \eta\epsilon \}_{1, \dots, N_r},$$

for which the binary classifier must predict the same class as for the boundary samples  $\mathcal{X}_b$ . Here,  $N_r \geq 0$  and  $\eta > 1.0$  control the number of samples and how far outside of the  $\epsilon$  ball they are located. Again, as for  $\mathcal{X}_b$ , we need to ensure that none of these samples get detected. By training the linear readout on  $\mathcal{X}_i$ ,  $\mathcal{X}_b$  and  $\mathcal{X}_r$  we guarantee that there exists at least one undetected adversarial sample within the  $\epsilon$ -ball around  $\mathbf{x}_c$ , and at least  $N_r$  samples outside the  $\epsilon$ -ball that are also undetected.

**Inverted Test** One potential issue with the normal test above, is that an attack might pass the test even though the attack completely ignores the detector. Indeed, many evaluations of detector defenses consider attacks that are oblivious to the presence of the detector [10]. Thus, an attack passing the test may not be sufficient to tell us that the attack is actually successfully targeting the detector.

To this end, we introduce a second *inverted test* that inverts the attack’s goal: Instead of finding adversarial samples that do not get detected, the goal is now to find an adversarial example that *is* detected. Since any detection defense that claims non-zero robustness must detect some adversarial examples, we can use these for constructing the set of boundary samples  $\mathcal{X}_b$ . Finally, we only need to negate the decision of the detection algorithm before proceeding exactly as for the previously described test.

Passing both the normal as well as the inverted test is a necessary condition for an adequate adversarial attack. In fact, this indicates that the attack is not agnostic to the detector but properly takes it into account. In contrast, passing only one of the tests indicates that only the classifier and not the detector is directly targeted.

## B. Experimental Details

All defenses investigated consider an  $\ell_\infty$  threat model. While the defense by Shan et al. [34] focuses on an  $\epsilon = 0.01$  bound, the rest uses the more common  $\epsilon = 8/255$  bound.

We evaluate the binarization test for 512 randomly chosen samples from the CIFAR-10 [18] test set.

For all attacks we set the gap between the boundary and inner points to  $\eta = 0.05$ , measured relatively to the used  $\epsilon$  value. We evaluated detector-based defenses using Algorithm 2, and use  $\xi = 1.75$ , measured in terms of  $\epsilon$ .

As outlined above in Section 4.3, we adjust the hardness of the test until the test produces conclusive results, i.e., the random attack success rate (R-ASR) is not too high. This leads to a parameter choice of  $N_{\text{inner}} = 999$  for all defenses but that of Zhang et al. [43] for which used

$N_{\text{inner}} = 9999$ . While we set the number of boundary samples to  $N_{\text{boundary}} = 10$  for Zhang et al. [44], we use set it to 1 for all other defenses. Also, we sample the boundary point(s) from the corners of the  $\ell_\infty$   $\epsilon$ -box, since this increases the test’s difficulty further.

Further, for adjusting the hardness of the test we adjust the bias of the linear classifier such that the distance between boundary sample and decision boundary measured in terms of the distance between boundary sample and closest inner sample is  $\kappa = 0.999$  (see Section 4.3).

We sample the inner samples uniformly from the  $\epsilon$  hypercube, and the boundary samples from the corners of the cube. We opted for this, since it increases the hardness of the test. Further, for calculating the R-ASR we samples both 200 points from the inner and 200 more from the corners of the space, as this significantly increased the R-ASR and, thus, gives a more realistic estimate of the test’s difficulty.

### C. Additional Results

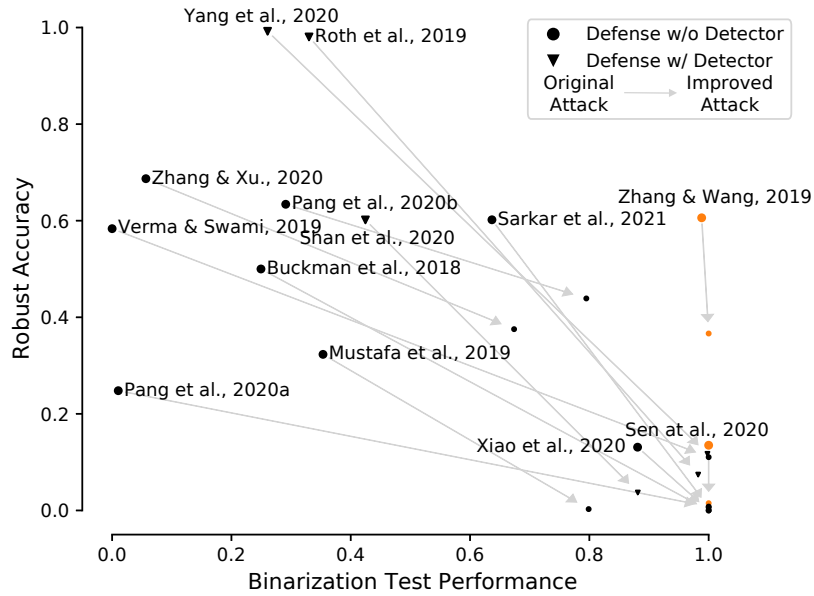


Figure 4. **Robust accuracy as a function of the test performance.** Thicker markers denote results for the attacks originally used by the defenses’ authors, while smaller ones correspond to that of adaptive attacks that broke the defense. The gray arrows between these points indicate how the scores change by using using a better suited attack. Orange points indicate false negatives/non-conclusive test results. Triangles denote defenses leveraging detection algorithms.