
RobustBench: a standardized adversarial robustness benchmark

Francesco Croce*
Univ. of Tübingen

Maksym Andriushchenko*
EPFL

Vikash Sehwal*
Princeton Univ.

Edoardo Debenedetti*
EPFL

Nicolas Flammarion
EPFL

Mung Chiang
Purdue Univ.

Prateek Mittal
Princeton Univ.

Matthias Hein
Univ. of Tübingen

Abstract

1 As a research community, we are still lacking a *systematic* understanding of the
2 progress on adversarial robustness which often makes it hard to identify the most
3 promising ideas in training robust models. A key challenge in benchmarking
4 robustness is that its evaluation is often error-prone leading to overestimation of
5 the true robustness of models. While adaptive attacks designed for a particular
6 defense are a potential solution, they have to be highly customized for particular
7 models, which makes it difficult to compare different methods. Our goal is to
8 instead establish a *standardized benchmark* of adversarial robustness, which as
9 accurately as possible reflects the robustness of the considered models within
10 a reasonable computational budget. To evaluate robustness of models for our
11 benchmark, we consider AutoAttack, an ensemble of white- and black-box attacks
12 which was recently shown in a large-scale study to improve almost all robustness
13 evaluations compared to the original publications. We also impose some restrictions
14 on the admitted models to rule out defenses that only make gradient-based attacks
15 ineffective without improving actual robustness. Our leaderboard, hosted at <http://robustbench.github.io/>,
16 contains evaluations of 90+ models and aims at reflecting the current state of the art on a set of well-defined tasks in ℓ_∞ - and ℓ_2 -
17 threat models and on common corruptions, with possible extensions in the future.
18 Additionally, we open-source the library <http://github.com/RobustBench/robustbench>
19 that provides unified access to 60+ robust models to facilitate their
20 downstream applications. Finally, based on the collected models, we analyze the
21 impact of robustness on the performance on distribution shifts, calibration, out-of-
22 distribution detection, fairness, privacy leakage, smoothness, and transferability.
23

24 1 Introduction

25 Since the finding that state-of-the-art deep learning models are vulnerable to small input perturbations
26 called *adversarial examples* [123], achieving adversarially robust models has become one of the most
27 studied topics in the machine learning community. The main difficulty of robustness evaluation is
28 that it is a computationally hard problem even for simple ℓ_p -bounded perturbations [64] and exact
29 approaches [126] do not scale to large enough models. There are already more than 3000 papers on
30 this topic [14], but it is often unclear which defenses against adversarial examples indeed improve
31 robustness and which only make the typically used attacks overestimate the actual robustness. There
32 is an important line of work on recommendations for how to perform adaptive attacks that are selected
33 specifically for a particular defense [4, 16, 129] which have in turn shown that several seemingly

*Equal contribution.

Rank	Method	Standard accuracy	Robust accuracy	Extra data	Architecture	Venue
1	Fixing Data Augmentation to Improve Adversarial Robustness	92.23%	66.56%	☑	WideResNet-70-16	arXiv, Mar 2021
2	Uncovering the Limits of Adversarial Training against Norm-Bounded Adversarial Examples	91.10%	65.87%	☑	WideResNet-70-16	arXiv, Oct 2020
3	Fixing Data Augmentation to Improve Adversarial Robustness	88.50%	64.58%	☒	WideResNet-106-16	arXiv, Mar 2021

Figure 1: The top-3 entries of our CIFAR-10 leaderboard hosted at <https://robustbench.github.io/> for the ℓ_∞ -perturbations of radius $\varepsilon_\infty = 8/255$.

robust defenses fail to be robust. However, recently Tramèr et al. [129] observe that although several recently published defenses have tried to perform adaptive evaluations, many of them could still be broken by new adaptive attacks. We observe that there are repeating patterns in many of these defenses that prevent standard attacks from succeeding. This motivates us to impose restrictions on the defenses we consider in our proposed benchmark, RobustBench, which aims at *standardized* adversarial robustness evaluation. Specifically, we rule out (1) classifiers which have zero gradients with respect to the input [12, 48], (2) randomized classifiers [147, 91], and (3) classifiers that contain an optimization loop in their predictions [108, 76]. Often, non-certified defenses that violate these three principles only make gradient-based attacks harder but do not substantially improve adversarial robustness [16]. We start from benchmarking robustness with respect to the ℓ_∞ - and ℓ_2 -threat models, since they are the most studied settings in the literature. We use the recent AutoAttack [26] as our current standard evaluation which is an ensemble of diverse parameter-free attacks (white- and black-box) that has shown for various datasets reliable performance over a large set of models that satisfy our restrictions. Moreover, we accept evaluations based on adaptive attacks whenever they can improve our standard evaluation. Additionally, we collect models robust against common image corruptions [53] as these represent another type of perturbations which should not modify the decision of a classifier although they are not produced in an adversarial way.

Contributions. We make the following contributions with our RobustBench benchmark:

- **Leaderboard** <https://robustbench.github.io/>: a website with the leaderboard (see Fig. 1) based on *more than 90* models where it is possible to track the progress and the current state of the art in adversarial robustness based on a standardized evaluation using AutoAttack (potentially complemented by adaptive attacks). The goal is to clearly identify the most successful ideas in training robust models to accelerate the progress in the field.
- **Model Zoo** <https://github.com/RobustBench/robustbench>: a collection of the most robust models that are easy to use for any downstream applications. For example, we expect that this will foster the development of better adversarial attacks by making it easier to perform evaluations on a large set of *more than 60* models.
- **Analysis**: based on the collected models from the Model Zoo, we provide an analysis of how robustness affects the performance on distribution shifts, calibration, out-of-distribution detection, fairness, privacy leakage, smoothness, and transferability. In particular, we find that robust models are significantly *underconfident* that leads to worse calibration, and that not all robust models have higher privacy leakage than standard models.

We believe that our standardized benchmark and accompanied collection of models will accelerate progress on multiple fronts in the area of adversarial robustness.

2 Background and related work

Adversarial perturbations. Let $\mathbf{x} \in \mathbb{R}^d$ be an input point and $y \in \{1, \dots, C\}$ be its correct label. For a classifier $f : \mathbb{R}^d \rightarrow \mathbb{R}^C$, we define a *successful adversarial perturbation* with respect to the perturbation set $\Delta \subseteq \mathbb{R}^d$ as a vector $\delta \in \mathbb{R}^d$ such that

$$\arg \max_{c \in \{1, \dots, C\}} f(\mathbf{x} + \delta)_c \neq y \quad \text{and} \quad \delta \in \Delta, \quad (1)$$

where typically the perturbation set Δ is chosen such that *all* points in $\mathbf{x} + \delta$ have y as their true label. This motivates a typical robustness measure called *robust accuracy*, which is the fraction of

74 datapoints on which the classifier f predicts the correct class for all possible perturbations from the
75 set Δ . Computing the exact robust accuracy is in general intractable and, when considering ℓ_p -balls
76 as Δ , NP-hard even for single-layer neural networks [64, 136]. In practice, an *upper bound* on the
77 robust accuracy is computed via some *adversarial attacks* which are mostly based on optimizing some
78 differentiable loss (e.g., cross entropy) using local search algorithms like projected gradient descent
79 (PGD) in order to find a successful adversarial perturbation. The tightness of the upper bound depends
80 on the effectiveness of the attack: unsuitable techniques or suboptimal parameters (e.g., the step size
81 and the number of iterations) can make the models appear more robust than they actually are [34, 86],
82 especially in the presence of phenomena like gradient obfuscation [4]. Certified methods [138, 44]
83 instead provide *lower bounds* on robust accuracy but often underestimate robustness significantly, in
84 particular if the certification was not part of the training process. Thus, we do not consider lower
85 bounds in our benchmark, and focus only on upper bounds which are typically much tighter [126].

86 **Threat models.** We focus on the fully white-box setting, i.e. the model f is assumed to be fully
87 known to the attacker. The threat model is defined by the set Δ of the allowed perturbations: the most
88 widely studied ones are the ℓ_p -perturbations, i.e. $\Delta_p = \{\delta \in \mathbb{R}^d, \|\delta\|_p \leq \varepsilon\}$, particularly for $p = \infty$
89 [123, 42, 79]. We rely on thresholds ε established in the literature which are chosen such that the true
90 label should stay the same for each in-distribution input within the perturbation set. We note that
91 robustness towards small ℓ_p -perturbations is a necessary but not sufficient notion of robustness which
92 has been criticized in the literature [41]. It is an active area of research to develop threat models
93 which are more aligned with the human perception such as spatial perturbations [39, 37], Wasserstein-
94 bounded perturbations [139, 57], perturbations of the image colors [72] or ℓ_p -perturbations in the
95 latent space of a neural network [73, 137]. However, despite the simplicity of the ℓ_p -perturbation
96 model, it has numerous interesting applications that go beyond security considerations [128, 106]
97 and span transfer learning [107, 132], interpretability [130, 65, 36], generalization [144, 158, 8],
98 robustness to unseen perturbations [62, 144, 73, 67], stabilization of GAN training [157]. Thus,
99 improvements in ℓ_p -robustness have the potential to improve many of these downstream applications.

100 **Common corruptions.** Unlike adversarial perturbations, common corruptions [53] try to mimic
101 modifications of the input images which can occur naturally: they are not imperceptible and evaluation
102 on them is done in the average case fashion, i.e. there is no attacker who aims at changing the
103 classifier’s decision. In this case, the robustness of a model is evaluated as classification accuracy on
104 the corrupted images, averaged over types and severities of corruptions.

105 **Related libraries and benchmarks.** There are many libraries that focus primarily on implemen-
106 tations of popular adversarial attacks such as FoolBox [100], Cleverhans [95], AdverTorch [31],
107 AdvBox [43], ART [89], SecML [83]. Some of them also provide implementations of several basic
108 defenses, but they do not include up-to-date state-of-the-art models.

109 The two challenges [71, 9] hosted at NeurIPS 2017 and 2018 aimed at finding the most robust models
110 for specific attacks, but they had a predefined deadline, so they could capture the best defenses only
111 at the time of the competition. Ling et al. [77] proposed DEEPSEC, a benchmark that tests many
112 combinations of attacks and defenses, but suffers from a few shortcomings as suggested by Carlini
113 [15], in particular: (1) reporting average-case performance over multiple attacks instead of worst-case
114 performance, (2) evaluating robustness in threat models different from the one used for training, (3)
115 using excessively large perturbations.

116 Recently, Dong et al. [33] have provided an evaluation of a few defenses (in particular, 3 for ℓ_∞ -
117 and 2 for ℓ_2 -norm on CIFAR-10) against multiple commonly used attacks. However, they did
118 not include some of the best performing defenses [55, 18, 46, 101] and attacks [45, 25], and in a
119 few cases, their evaluation suggests robustness higher than what was reported in the original papers.
120 Moreover, they do not impose any restrictions on the models they accept to the benchmark. RobustML
121 (<https://www.robust-ml.org/>) aims at collecting robustness claims for defenses together with
122 external evaluations. Their format does not assume running any baseline attack, so it relies entirely
123 on evaluations submitted by the community, which however do not occur often enough. Thus even
124 though RobustML has been a valuable contribution to the community, now it does not provide a
125 comprehensive overview of the recent state of the art in adversarial robustness.

126 Finally, it has become common practice to test new attacks wrt ℓ_∞ on the publicly available models
127 from Madry et al. [79] and Zhang et al. [154], since those represent widely accepted defenses which
128 have stood many thorough evaluations. However, having only two models per dataset (MNIST and

129 CIFAR-10) does not constitute a sufficiently large testbed, and, because of the repetitive evaluations,
130 some attacks may already overfit to those defenses.

131 **What is different in RobustBench.** Learning from these previous attempts, RobustBench presents
132 a few different features compared to the aforementioned benchmarks: (1) a baseline worst-case
133 evaluation with an ensemble of *strong, standardized* attacks [26] which includes both white- and
134 black-box attacks that can be *optionally* extended by adaptive evaluations, (2) clearly defined threat
135 models that correspond to the ones used during training for submitted defenses, (3) evaluation of not
136 only standard defenses [79] but also of more recent improvements such as [18, 46, 101], (4) the Model
137 Zoo that provides convenient access to the 60+ most robust models from the literature which can be
138 used for downstream tasks and facilitate the development of new standardized attacks. Moreover,
139 RobustBench is designed as an *open-ended* benchmark that keeps an up-to-date leaderboard, and
140 we welcome contributions of new defenses and evaluations of adaptive attacks for particular models.

141 3 Description of RobustBench

142 In this section, we start by providing a detailed layout of our proposed leaderboard for ℓ_∞ , ℓ_2 , and
143 the common corruptions threat models. Next, we present the Model Zoo, which provides unified
144 access to most networks from our leaderboards.

145 3.1 Leaderboard

146 **Restrictions.** We argue that benchmarking adversarial robustness in a standardized way requires
147 some restrictions on the type of considered models. The goal of these restrictions is to prevent
148 submissions of defenses that cause some standard attacks to fail without actually improving robustness.
149 Specifically, we consider only classifiers $f : \mathbb{R}^d \rightarrow \mathbb{R}^C$ that

- 150 • have in general *non-zero gradients* with respect to the inputs. Models with zero gradients,
151 e.g., that rely on quantization of inputs [12, 48], make gradient-based methods ineffective
152 thus requiring zeroth-order attacks, which do not perform as well as gradient-based attacks.
153 Alternatively, specific adaptive evaluations, e.g. with Backward Pass Differentiable Approx-
154 imation [4], can be used which, however, can hardly be standardized. Moreover, we are not
155 aware of existing defenses solely based on having zero gradients for large parts of the input
156 space which would achieve competitive robustness.
- 157 • have a *fully deterministic forward pass*. To evaluate defenses with stochastic components,
158 it is a common practice to combine standard gradient-based attacks with Expectation over
159 Transformations [4]. While often effective, it might be not sufficient, as shown by Tramèr
160 et al. [129]. Moreover, the classification decision of randomized models may vary over
161 different runs for the same input, hence even the definition of robust accuracy differs from
162 that of deterministic networks. We also note that randomization *can* be useful for improving
163 robustness and deriving robustness certificates [74, 23], but it also introduces variance in the
164 gradient estimators (both white- and black-box) which can make attacks much less effective.
- 165 • do not have an *optimization loop* in the forward pass. This makes backpropagation through
166 the classifier very difficult or extremely expensive. Usually, such defenses [108, 76] need to
167 be evaluated adaptively with attacks considering jointly the loss of the inner loop and the
168 standard classification task.

169 Some of these restrictions were also discussed by [11] for the warm-up phase of their challenge. We
170 refer the reader to Appendix E therein for an illustrative example of a trivial defense that bypasses
171 gradient-based and some of the black-box attacks they consider.

172 **Overall setup.** We set up leaderboards for the ℓ_∞ , ℓ_2 and common corruption threat models on
173 CIFAR-10 and CIFAR-100 [69] datasets (see Table 1 for details). We use the fixed budgets of
174 $\varepsilon_\infty = 8/255$ and $\varepsilon_2 = 0.5$ for the ℓ_∞ and ℓ_2 leaderboards. Most of the models shown there are taken
175 from papers published at top-tier machine learning and computer vision conferences as shown in
176 Fig. 2 (left). For each entry we report the reference to the original paper, standard and robust accuracy
177 under the specific threat model (see the next paragraph for details), network architecture, venue
178 where the paper appeared and possibly notes regarding the model. We also highlight when extra data
179 (usually, the dataset introduced by Carmon et al. [18]) is used since it gives a clear advantage for both

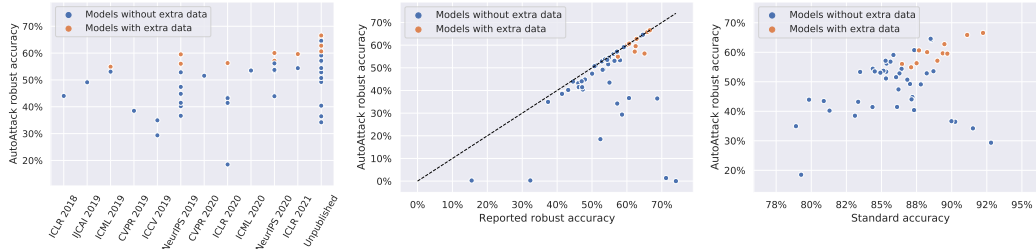


Figure 2: Visualization of the robustness and accuracy of 54 CIFAR-10 models from the RobustBench ℓ_∞ -leaderboard. Robustness is evaluated using ℓ_∞ -perturbations with $\varepsilon_\infty = 8/255$.

180 clean and robust accuracy. Moreover, the leaderboard allows to search the entries by their metadata
 181 (such as title, architecture, venue) which can be useful to compare different methods that use the
 182 same architecture or to search for papers published at some recent conference.

183 **Evaluation of defenses.** The evaluation of robust accuracy on common corruptions [53] involves
 184 simply computing the average accuracy on corrupted images over different corruption types and
 185 severity levels.¹ To evaluate robustness of ℓ_∞ and ℓ_2 defenses, we currently use AutoAttack [26].
 186 It is an ensemble of four attacks: a variation of PGD attack with automatically adjusted step sizes,
 187 with (1) the cross entropy loss and (2) the difference of logits ratio loss, which is a rescaling-invariant
 188 margin-based loss function, (3) the targeted version of the FAB attack [25], which minimizes the
 189 ℓ_p -norm of the perturbations, and (4) the black-box Square Attack [3]. We choose AutoAttack as it
 190 includes both black-box and white-box attacks, does not require hyperparameter tuning (in particular,
 191 the step size), and consistently improves the results reported in the original papers for almost all the
 192 models (see Fig. 2 (middle)). If in the future some new standardized and parameter-free attack is
 193 shown to consistently outperform AutoAttack on a wide set of models given a similar computational
 194 cost, we will adopt it as standard evaluation. In order to verify the reproducibility of the results, we
 195 perform the standardized evaluation independently of the authors of the submitted models. Below we
 196 show an example of how one can use our library to easily benchmark a model (either external one or
 197 taken from the Model Zoo):

```
from robustbench.eval import benchmark
clean_acc, robust_acc = benchmark(model, dataset='cifar10', threat_model='Linf')
```

198 Moreover, in Appendix C we also show the variability of the robust accuracy given by AutoAttack over
 199 random seeds and report its runtime for a few models from different threat models. We also accept
 200 evaluations of the individual models on the leaderboard based on adaptive or external attacks to reflect
 201 the best available upper bound on the true robust accuracy. For example, Goyal et al. [46] and Rebuffi
 202 et al. [101] evaluate their models with a hybrid of AutoAttack and MultiTargeted attack [45], that in
 203 some cases report slightly lower robust accuracy than AutoAttack alone. We reflect all such additional
 204 evaluations in our leaderboard. [The submission of adaptive evaluations is facilitated by a pre-formatted
 205 issue template in our repository https://github.com/RobustBench/robustbench.](https://github.com/RobustBench/robustbench)

206 **Adding new defenses.** We believe that the leaderboard is only useful if it reflects the latest advances
 207 in the field, so it needs to be constantly updated with new defenses. We intend to include evaluations
 208 of new techniques and we welcome contributions from the community which can help to keep the
 209 benchmark up-to-date. We require new entries to (1) satisfy the three restrictions stated above, (2)
 210 to be accompanied by a publicly available paper (e.g., an arXiv preprint) describing the technique
 211 used to achieve the reported results, and (3) share the model checkpoints (not necessarily publicly).
 212 We also allow *temporarily* adding entries without providing checkpoints given that the authors
 213 evaluate their models with AutoAttack. However, we will mark such evaluations as *unverified*, and to
 214 encourage reproducibility, we reserve the right to remove an entry later on if the corresponding model
 215 checkpoint is not provided. It is possible to add a new defense to the leaderboard and (optionally)
 216 the Model Zoo by opening an issue with a predefined template in our repository [https://github.
 217 com/RobustBench/robustbench](https://github.com/RobustBench/robustbench), where more details about new additions can be found.

¹A breakdown over corruptions and severities is also available, e.g. for CIFAR-10 models see: https://github.com/RobustBench/robustbench/blob/master/model_info/cifar10/corruptions/unaggregated_results.csv

Table 1: The total number of models in the Model Zoo and leaderboards per dataset and threat model.

Threat model	CIFAR-10		CIFAR-100	
	Model Zoo	Leaderboard	Model Zoo	Leaderboard
ℓ_∞ with $\varepsilon_\infty = 8/255$	33	55	12	12
ℓ_2 with $\varepsilon_2 = 0.5$	14	14	-	-
Common corruptions [53]	7	12	2	4

218 3.2 Model Zoo

219 We collect the checkpoints of many networks from the leaderboard in a single repository hosted at
 220 <https://github.com/RobustBench/robustbench> after obtaining the permission of the authors
 221 (see Appendix A for the information on the licenses). The goal of this repository, the Model Zoo, is to
 222 make the usage of robust models as simple as possible to facilitate various downstream applications
 223 and analyses of general trends in the field. In fact, even when the checkpoints of the proposed method
 224 are made available by the authors, it is often time-consuming and not straightforward to integrate them
 225 in the same framework because of many factors such as small variations in the architectures, custom
 226 input normalizations, etc. For simplicity of implementation, at the moment we include only models
 227 implemented in PyTorch [96]. Below we illustrate how a model can be automatically downloaded
 228 and loaded via its identifier and threat model within two lines of code:

```

 229 from robustbench.utils import load_model
 230 model = load_model(model_name='Carmon2019Unlabeled',
 231                   dataset='cifar10', threat_model='Linf')
```

229 At the moment, all models (see Table 1 and Appendix E for details) are variations of ResNet [50] and
 230 WideResNet architectures [150] of different depth and width. We include the most robust models, e.g.
 231 those from Rebuffi et al. [101], but there are also defenses which pursue additional goals alongside
 232 adversarial robustness at the fixed threshold we use: e.g., Sehwag et al. [112] consider networks
 233 which are robust and compact, Wong et al. [140] focus on computationally efficient adversarial
 234 training, Ding et al. [32] aim at input-adaptive robustness as opposed to robustness within a single
 235 ℓ_p -radius. All these factors have to be taken into account when comparing different techniques, as
 236 they have a strong influence on the final performance.

237 **A testbed for new attacks.** Another important use case of the Model Zoo is to simplify comparisons
 238 between different adversarial attacks on a wide range of models. First, the leaderboard already serves
 239 as a strong baseline for new attacks. Second, as mentioned above, new attacks are often evaluated on
 240 the models from Madry et al. [79] and Zhang et al. [154], but this may not provide a representative
 241 picture of their effectiveness. For example, currently the difference in robust accuracy between the
 242 first and second-best attacks in the CIFAR-10 leaderboard of Madry et al. [79] is only 0.03%, and
 243 between the second and third is 0.04%. Thus, we believe that a more thorough comparison should
 244 involve multiple models to prevent overfitting of the attack to one or two standard robust defenses.

245 4 Analysis

246 With unified access to multiple models from the Model Zoo, one can easily compute various perfor-
 247 mance metrics to see general trends. In the following we analyze various aspects of robust classifiers,
 248 reporting results mostly for ℓ_∞ -robust models on CIFAR-10 while the results for other threat models
 249 and datasets can be found in Appendix D.

250 **Progress on adversarial defenses.** In Fig. 2, we plot a breakdown over conferences, the amount
 251 of robustness overestimation reported in the original papers, and we also visualize the robustness-
 252 accuracy trade-off for the ℓ_∞ -models from the Model Zoo. First, we observe that for multiple
 253 *published* defenses, the reported robust accuracy is highly overestimated. We also find that the use of
 254 extra data is able to alleviate the robustness-accuracy trade-off as suggested in previous works [98].
 255 However, so far all models with high robustness to perturbations of ℓ_∞ -norm up to $\varepsilon = 8/255$ still
 256 suffer from noticeable degradation in clean accuracy compared to standardly trained models. Finally,
 257 it is interesting to note that the best entries of the ℓ_p -leaderboards are still variants of PGD adversarial
 258 training [79, 154] but with various enhancements (extra data, early stopping, weight averaging).

259 **Performance across various distribution shifts.** Here we test the performance of the models from
 260 the Model Zoo on different distribution shifts ranging from common image corruptions (CIFAR-10-C,
 261 [53]) to dataset resampling bias (CIFAR-10.1, [102]) and image source shift (CINIC-10, [29]). For
 262 each of these datasets, we measure standard accuracy, and Fig. 3 shows that improvement in the
 263 robust accuracy (which often comes with an improvement in standard accuracy) on CIFAR-10 also
 264 correlates with an improvement in standard accuracy across distributional shifts. On CIFAR-10-C,
 265 we observe that robust models (particularly with respect to the ℓ_2 -norm) tend to give a significant
 266 improvement which agrees with the findings from the previous literature [40]. Concurrently with our
 267 work, Taori et al. [125] also study the robustness to different distribution shifts of many models trained
 268 on ImageNet, including some ℓ_p -robust models. Our conclusions qualitatively agree with theirs, and
 269 we hope that our collected set of models will help to provide a more complete picture. Moreover,
 270 we measure robust accuracy, in the same threat model used on CIFAR-10, using AutoAttack [26]
 271 (see Fig. 10 in Appendix D), and notice how ℓ_p adversarial robustness generalizes across different
 datasets, and a clear positive correlation between robust accuracy on CIFAR-10 and its variations.

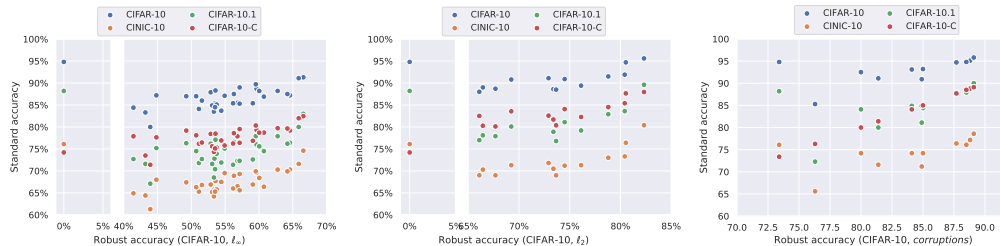


Figure 3: Standard accuracy of classifiers trained against ℓ_∞ (left), ℓ_2 (middle), and common corruption (right) threat model respectively, from our Model Zoo on various distribution shifts.

272

273 **Calibration.** A classifier is *calibrated* if its predicted probabilities correctly reflect the actual accuracy
 274 [47]. In the context of adversarial training, calibration was considered in Hendrycks et al. [56] who
 275 focus on improving accuracy on common corruptions and in Augustin et al. [6] who focus mostly on
 276 preventing overconfident predictions on out-of-distribution inputs. We instead focus on *in-distribution*
 277 calibration, and in Fig. 4 plot the expected calibration error (ECE) without and with temperature
 278 rescaling [49] to minimize the ECE (which is a simple but effective post-hoc calibration method,
 279 see Appendix D for details) together with the optimal temperature for a large set of ℓ_∞ models.
 280 We observe that most of the ℓ_∞ robust models are significantly *underconfident* since the optimal
 281 calibration temperature is less than one for most models. The only two models in Fig. 4 which are
 282 *overconfident* are the standard model and the model of Ding et al. [32] that aims to maximize the
 283 margin. We see that temperature rescaling is even more important for robust models since without
 284 any rescaling the ECE is as high as 70% for the model of Pang et al. [92] (and 21% on average)
 285 compared to 4% for the standard model. Temperature rescaling significantly reduces the ECE gap
 286 between robust and standard models but it does not fix the problem completely which suggests that it
 287 is worth incorporating calibration techniques also during training of robust models. For ℓ_2 robust
 288 models, the models can be on the contrary *more calibrated* by default, although the improvement
 vanishes if temperature rescaling is applied (see Appendix D).

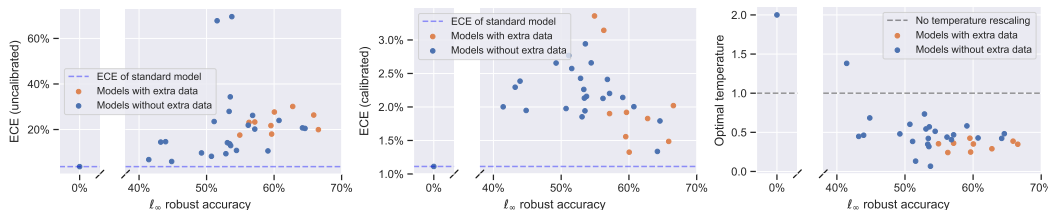


Figure 4: Expected calibration error (ECE) before (left) and after (middle) temperature rescaling, and the optimal rescaling temperature (right) for the ℓ_∞ -robust models.

289

290 **Out-of-distribution detection.** Ideally, a classifier should exhibit high uncertainty in its predictions
 291 when evaluated on *out-of-distribution* (OOD) inputs. One of the most straightforward ways to extract
 292 this uncertainty information is to use some threshold on the predicted confidence where OOD inputs

293 are expected to have low confidence from the model [54]. An emerging line of research aims at
 294 developing OOD detection methods in conjunction with adversarial robustness [52, 110, 6]. In
 295 particular, Song et al. [122] demonstrated that adversarial training [79] leads to degradation in the
 296 robustness against OOD data. We further test this observation on all ℓ_∞ -models trained on CIFAR-10
 297 from the Model Zoo on three OOD datasets: CIFAR-100 [69], SVHN [88], and Describable Textures
 298 Dataset [22]. We use the area under the ROC curve (AUROC) to measure the success in the detection
 299 of OOD data, and show the results in Fig. 5. With ℓ_∞ robust models, we find that compared to
 300 standard training, various robust training methods indeed lead to degradation of the OOD detection
 301 quality. While extra data in standard training can improve robustness against OOD inputs, it fails
 302 to provide similar improvements with robust training. We further find that ℓ_2 robust models have in
 303 general comparable OOD detection performance to standard models (see Fig. 12 in Appendix), while
 304 the model of Augustin et al. [6] achieves even better performance since their approach explicitly
 optimizes both robust accuracy and worst-case OOD detection performance.

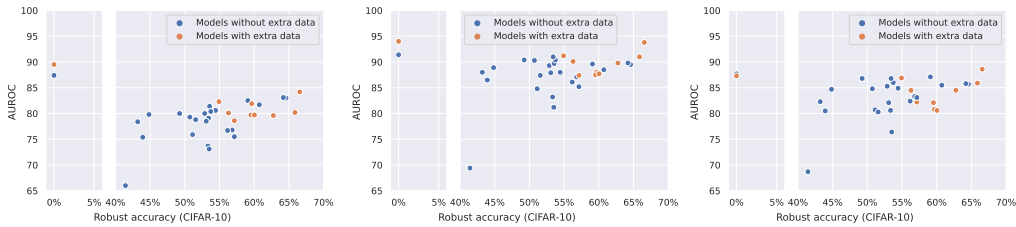


Figure 5: Visualization of the OOD detection quality (higher AUROC is better) for the ℓ_∞ -robust models trained on CIFAR-10 on three OOD datasets: CIFAR-100 (left), SVHN (middle), Describable Textures (right). We detect OOD inputs based on the maximum predicted confidence [54].

305

306 **Fairness in robustness.** Recent works [7, 146] have noticed that robust training [79, 154] can lead
 307 to models whose performance varies significantly across subgroups, e.g. defined by classes. We will
 308 refer to this performance difference as *fairness*, and here we study the influence of robust training
 309 methods on fairness. In Fig. 6 we show the breakdown of standard and robust accuracy for the ℓ_∞
 310 robust models, where one can see how the achieved robustness largely varies over classes. While in
 311 general the classwise standard and robust accuracy correlate well, the class “deer” in ℓ_∞ -threat model
 312 suffers a significant degradation, unlike what happens for ℓ_2 (see Appendix D), which might indicate
 313 that the features of such class are particularly sensitive to ℓ_∞ -bounded attacks. Moreover, we measure
 314 fairness with the relative standard deviation (RSD), defined as the standard deviation divided over the
 315 average, of robust accuracy over classes for which lower values mean more uniform distribution and
 316 higher robustness. We observe that better robust accuracy generally leads to lower RSD values which
 317 implies that the disparity among classes is reduced. However, some training techniques like MART
 318 [135] can noticeably increase the RSD and thus *increase the disparity* compared to other methods
 319 which achieve similar robustness (around 57%).

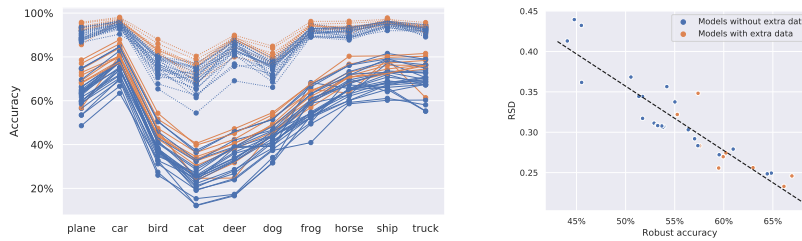


Figure 6: Fairness of ℓ_∞ -robust models. **Left:** classwise standard (dotted lines) and robust (solid) accuracy. **Right:** relative standard deviation (RSD) of robust accuracy over classes vs its average.

320 **Privacy leakage.** Deep neural networks are prone to memorizing training data [117, 17]. Recent
 321 works have highlighted that robust training exacerbates this problem [121]. Here we benchmark
 322 privacy leakage of training data across robust networks (Fig. 7). We calculate membership inference
 323 accuracy using output confidence of adversarial images from the training and test sets (see Appendix D
 324 for more details). It measures how accurately we can infer whether a sample was present in the
 325 training dataset. Our analysis reveals mixed trends. First, our results show that not all robust models
 326 have a significantly higher privacy leakage than a standard model. We find that the inference accuracy

327 across robust models has a large variation, where some models even have lower privacy leakage than
 328 a standard model. It also does not have a strong correlation with the robust accuracy. In contrast, it
 329 is largely determined by the generalization gap, as using classification confidence information does
 330 not lead to a much higher inference accuracy than the baseline determined by the generalization
 331 gap (as shown in Fig. 7 (right)). Thus one can expect lower privacy leakage in robust networks as
 332 multiple previous works have explicitly aimed to reduce the generalization gap in robust training using
 333 techniques such as early stopping [103, 154, 46]. It further suggests that reducing the generalization
 334 gap in robust networks can further reduce privacy leakage.



Figure 7: Privacy leakage of ℓ_∞ -robust models. We measure privacy leakage of training data in robust networks and compare it with robust accuracy (left) and generalization gap (right).

335 **Extra experiments.** In Appendix D, we show extra experiments related to the points analyzed above
 336 and describe some of the implementation details. Also, we study how adversarial perturbations
 337 transfer between different models. We find that adversarial examples strongly transfer from robust
 338 to robust, non-robust to robust, and non-robust to non-robust networks. However, we observe poor
 339 transferability of adversarial examples from robust to non-robust networks. Moreover, since prior
 340 works [51, 148] connected higher smoothness with better robustness, we analyze the smoothness
 341 of the models both at intermediate and output layers. This confirms that, for a fixed architecture,
 342 standard training yields classifiers that are significantly less smooth than robust ones. This illustrates
 343 that one can use the collected models to study the *internal* properties of robust networks as well.

344 5 Outlook

345 **Conclusions.** We believe that a *standardized* benchmark with clearly defined threat models, restric-
 346 tions on submitted models, and tight upper bounds on robust accuracy can be useful to show which
 347 ideas in training robust models are the most successful. Recent works have already referred to our
 348 leaderboards [68, 149, 80, 124, 145], in particular as reflecting the current state of the art [101, 75, 94],
 349 and used the networks of our Model Zoo to test new adversarial attacks [83, 105, 38, 109], to evaluate
 350 test-time defenses [133] or to evaluate perceptual distances derived from them [61]. Additionally,
 351 we have shown that unified access to a *large* and *up-to-date* set of robust models can be useful to
 352 analyze multiple aspects related to robustness. First, one can easily analyze the progress of adversarial
 353 defenses over time including the amount of robustness overestimation and the robustness-accuracy
 354 tradeoff. Second, one can conveniently study the impact of robustness on other performance metrics
 355 such as accuracy under distribution shifts, calibration, out-of-distribution detection, fairness, privacy
 356 leakage, smoothness, and transferability. Overall, we think that the community has to develop a better
 357 understanding of how different types of robustness affect other aspects of the model performance and
 358 RobustBench can help to achieve this goal.

359 **Broader impact.** In our work, we do not only perform a standardized benchmarking of adversarial
 360 robustness but also analyze multiple other properties of robust models such as calibration, privacy
 361 leakage, fairness, etc. Such analyses are important, in our opinion, since they allow us to assess the
 362 broader impact of improving robustness on other crucial performance metrics of neural networks. Ad-
 363 ditionally, in motivating higher robustness against adversarial examples, our work leaves an unwanted
 364 side effect on tasks where adversarial attacks can actually be used for beneficial purposes [115, 99].
 365 Finally, we note that a good performance on our benchmark does not guarantee the safety of the
 366 benchmarked model in a real-world deployment since ℓ_p - and corruption robustness may not be
 367 necessarily a realistic threat model (although it is a insightful problem to work on) and the real-world
 368 robustness is likely to require more domain-specific threat models.

369 **Future plans.** Our intention in the future is to keep the current leaderboards up-to-date (see the
 370 maintenance plan in Appendix B) and add new leaderboards for other datasets (in particular, for
 371 ImageNet [30]) and other threat models which become widely accepted in the community. We see
 372 as potential candidates (1) sparse perturbations, e.g. bounded by ℓ_0 , ℓ_1 -norm or adversarial patches
 373 [10, 24, 84, 27], (2) multiple ℓ_p -norm perturbations [127, 81], (3) adversarially optimized common
 374 corruptions [62, 63], (4) a broad set of perturbations unseen during training [73]. Another possible
 375 direction of development of the benchmark is including defenses based on some form of test-time
 376 adaptation [116, 133], which do not fulfill the third restriction (no optimization loop). However, since
 377 those are showing promising results and drawing attention from the community, one can introduce a
 378 separate leaderboard with specific rules and evaluation protocol for them.

379 References

- 380 [1] M. Alfara, J. C. Perez, A. Bibi, A. Thabet, P. Arbelaez, and B. Ghanem. Clustr: Clustering training for
 381 robustness. *arXiv*, 2020.
- 382 [2] M. Andriushchenko and N. Flammarion. Understanding and improving fast adversarial training. *NeurIPS*,
 383 2020.
- 384 [3] M. Andriushchenko, F. Croce, N. Flammarion, and M. Hein. Square attack: a query-efficient black-box
 385 adversarial attack via random search. In *ECCV*, 2020.
- 386 [4] A. Athalye, N. Carlini, and D. Wagner. Obfuscated gradients give a false sense of security: Circumventing
 387 defenses to adversarial examples. In *ICML*, 2018.
- 388 [5] M. Atzmon, N. Haim, L. Yariv, O. Israelov, H. Maron, and Y. Lipman. Controlling neural level sets.
 389 *NeurIPS*, 2019.
- 390 [6] M. Augustin, A. Meinke, and M. Hein. Adversarial robustness on in- and out-distribution improves
 391 explainability. *ECCV*, 2020.
- 392 [7] P. Benz, C. Zhang, A. Karjauv, and I. S. Kweon. Robustness may be at odds with fairness: An empirical
 393 study on class-wise accuracy. *arXiv preprint arXiv:2010.13365*, 2020.
- 394 [8] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao. Yolov4: Optimal speed and accuracy of object detection.
 395 *arXiv preprint arXiv:2004.10934*, 2020.
- 396 [9] W. Brendel, J. Rauber, A. Kurakin, N. Papernot, B. Velicki, M. Salathé, S. P. Mohanty, and M. Bethge.
 397 Adversarial vision challenge. In *NeurIPS Competition Track*, 2018.
- 398 [10] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer. Adversarial patch. In *NeurIPS 2017 Workshop*
 399 *on Machine Learning and Computer Security*, 2017.
- 400 [11] T. B. Brown, N. Carlini, C. Zhang, C. Olsson, P. Christiano, and I. Goodfellow. Unrestricted adversarial
 401 examples. *arXiv preprint arXiv:1809.08352*, 2018.
- 402 [12] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow. Thermometer encoding: One hot way to resist
 403 adversarial examples. In *ICLR*, 2018.
- 404 [13] D. A. Calian, F. Stimberg, O. Wiles, S.-A. Rebuffi, A. Gyorgy, T. Mann, and S. Goyal. Defending against
 405 image corruptions through adversarial augmentations. *arXiv*, 2021.
- 406 [14] N. Carlini. A complete list of all (arxiv) adversarial example papers. <https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>. Accessed: 2021-06-08.
- 407 [15] N. Carlini. A critique of the deepsec platform for security analysis of deep learning models. *arXiv*
 408 *preprint arXiv:1905.07112*, 2019.
- 410 [16] N. Carlini, A. Athalye, N. Papernot, W. Brendel, J. Rauber, D. Tsipras, I. Goodfellow, A. Madry, and
 411 A. Kurakin. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.
- 412 [17] N. Carlini, C. Liu, Ú. Erlingsson, J. Kos, and D. Song. The secret sharer: Evaluating and testing
 413 unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX Security*
 414 *19)*, pages 267–284, 2019.
- 415 [18] Y. Carmon, A. Raghunathan, L. Schmidt, P. Liang, and J. C. Duchi. Unlabeled data improves adversarial
 416 robustness. *NeurIPS*, 2019.
- 417 [19] A. Chan, Y. Tay, Y. S. Ong, and J. Fu. Jacobian adversarially regularized networks for robustness. *ICLR*,
 418 2020.
- 419 [20] J. Chen, Y. Cheng, Z. Gan, Q. Gu, and J. Liu. Efficient robust training via backward smoothing. *arXiv*,
 420 2020.
- 421 [21] T. Chen, S. Liu, S. Chang, Y. Cheng, L. Amini, and Z. Wang. Adversarial robustness: From self-supervised
 422 pre-training to fine-tuning. In *CVPR*, 2020.

- 423 [22] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *CVPR*,
424 2014.
- 425 [23] J. M. Cohen, E. Rosenfeld, and J. Z. Kolter. Certified adversarial robustness via randomized smoothing.
426 In *ICML*, 2019.
- 427 [24] F. Croce and M. Hein. Sparse and imperceptible adversarial attacks. In *ICCV*, 2019.
- 428 [25] F. Croce and M. Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In
429 *ICML*, 2020.
- 430 [26] F. Croce and M. Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse
431 parameter-free attacks. In *ICML*, 2020.
- 432 [27] F. Croce, M. Andriushchenko, N. D. Singh, N. Flammarion, and M. Hein. Sparse-rs: a versatile framework
433 for query-efficient sparse black-box adversarial attacks. In *ECCV Workshop on Adversarial Robustness in*
434 *the Real World*, 2020.
- 435 [28] J. Cui, S. Liu, L. Wang, and J. Jia. Learnable boundary guided adversarial training. *arXiv*, 2020.
- 436 [29] L. N. Darlow, E. J. Crowley, A. Antoniou, and A. J. Storkey. Cinic-10 is not imagenet or cifar-10. *arXiv*
437 *preprint arXiv:1810.03505*, 2018.
- 438 [30] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image
439 database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee,
440 2009.
- 441 [31] G. W. Ding, L. Wang, and X. Jin. AdverTorch v0.1: An adversarial robustness toolbox based on pytorch.
442 *arXiv preprint arXiv:1902.07623*, 2019.
- 443 [32] G. W. Ding, Y. Sharma, K. Y. C. Lui, and R. Huang. Mma training: Direct input space margin
444 maximization through adversarial training. In *ICLR*, 2020.
- 445 [33] Y. Dong, Q.-A. Fu, X. Yang, T. Pang, H. Su, Z. Xiao, and J. Zhu. Benchmarking adversarial robustness
446 on image classification. In *CVPR*, 2020.
- 447 [34] L. Engstrom, A. Ilyas, and A. Athalye. Evaluating and understanding the robustness of adversarial logit
448 pairing. *NeurIPS 2018 Workshop on Security in Machine Learning*, 2018.
- 449 [35] L. Engstrom, A. Ilyas, H. Salman, S. Santurkar, and D. Tsipras. Robustness (python library), 2019. URL
450 <https://github.com/MadryLab/robustness>.
- 451 [36] L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, B. Tran, and A. Madry. Adversarial robustness as a prior
452 for learned representations. *arXiv preprint arXiv:1906.00945*, 2019.
- 453 [37] L. Engstrom, B. Tran, D. Tsipras, L. Schmidt, and A. Madry. Exploring the landscape of spatial robustness.
454 In *ICML*, 2019.
- 455 [38] F. Faghri, C. Vasconcelos, D. J. Fleet, F. Pedregosa, and N. L. Roux. Bridging the gap between adversarial
456 robustness and optimization bias. *arXiv preprint arXiv:2102.08868*, 2021.
- 457 [39] A. Fawzi and P. Frossard. Manitest: Are classifiers really invariant? In *BMVC*, 2015.
- 458 [40] N. Ford, J. Gilmer, N. Carlini, and D. Cubuk. Adversarial examples are a natural consequence of test
459 error in noise. In *ICML*, 2019.
- 460 [41] J. Gilmer, R. P. Adams, I. Goodfellow, D. Andersen, and G. E. Dahl. Motivating the rules of the game for
461 adversarial example research. *arXiv preprint arXiv:1807.06732*, 2018.
- 462 [42] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *ICLR*,
463 2015.
- 464 [43] D. Goodman, H. Xin, W. Yang, W. Yuesheng, X. Junfeng, and Z. Huan. Advbox: a toolbox to generate
465 adversarial examples that fool neural networks. *arXiv preprint arXiv:2001.05574*, 2020.
- 466 [44] S. Gowal, K. D. Dvijotham, R. Stanforth, R. Bunel, C. Qin, J. Uesato, R. Arandjelovic, T. Mann, and
467 P. Kohli. Scalable verified training for provably robust image classification. In *ICCV*, 2019.
- 468 [45] S. Gowal, J. Uesato, C. Qin, P.-S. Huang, T. Mann, and P. Kohli. An alternative surrogate loss for
469 pgd-based adversarial testing. *arXiv preprint arXiv:1910.09338*, 2019.
- 470 [46] S. Gowal, C. Qin, J. Uesato, T. Mann, and P. Kohli. Uncovering the limits of adversarial training against
471 norm-bounded adversarial examples. *arXiv*, 2020.
- 472 [47] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In
473 *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.
- 474 [48] C. Guo, M. Rana, M. Cisse, and L. van der Maaten. Countering adversarial images using input transfor-
475 mations. In *ICLR*, 2018.
- 476 [49] C. Guo, J. R. Gardner, Y. You, A. G. Wilson, and K. Q. Weinberger. Simple black-box adversarial attacks.
477 In *ICML*, 2019.

- 478 [50] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- 479 [51] M. Hein and M. Andriushchenko. Formal guarantees on the robustness of a classifier against adversarial
480 manipulation. In *NeurIPS*, 2017.
- 481 [52] M. Hein, M. Andriushchenko, and J. Bitterwolf. Why relu networks yield high-confidence predictions far
482 away from the training data and how to mitigate the problem. In *CVPR*, 2019.
- 483 [53] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions and
484 perturbations. In *ICLR*, 2019.
- 485 [54] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in
486 neural networks. In *ICLR*, 2017.
- 487 [55] D. Hendrycks, K. Lee, and M. Mazeika. Using pre-training can improve model robustness and uncertainty.
488 In *ICML*, 2019.
- 489 [56] D. Hendrycks, N. Mu, E. D. Cubuk, B. Zoph, J. Gilmer, and B. Lakshminarayanan. Augmix: A simple
490 data processing method to improve robustness and uncertainty. *ICLR*, 2020.
- 491 [57] J. E. Hu, A. Swaminathan, H. Salman, and G. Yang. Improved image wasserstein attacks and defenses.
492 *ICLR Workshop: Towards Trustworthy ML: Rethinking Security and Privacy for ML*, 2020.
- 493 [58] L. Huang, C. Zhang, and H. Zhang. Self-adaptive training: beyond empirical risk minimization. *NeurIPS*,
494 2020.
- 495 [59] Y. Jang, T. Zhao, S. Hong, and H. Lee. Adversarial defense via learning to generate diverse attacks. *ICCV*,
496 2019.
- 497 [60] C. Jin and M. Rinard. Manifold regularization for adversarial robustness. *arXiv*, 2020.
- 498 [61] A. Ju. Generative models as a robust alternative for image classification: Progress and challenges. *PhD*
499 *thesis, UC Berkeley*, 2021.
- 500 [62] D. Kang, Y. Sun, T. Brown, D. Hendrycks, and J. Steinhardt. Transfer of adversarial robustness between
501 perturbation types. *arXiv preprint arXiv:1905.01034*, 2019.
- 502 [63] D. Kang, Y. Sun, D. Hendrycks, T. Brown, and J. Steinhardt. Testing robustness against unforeseen
503 adversaries. *arXiv preprint arXiv:1908.08016*, 2019.
- 504 [64] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer. Reluplex: an efficient smt solver for
505 verifying deep neural networks. In *ICCAV*, 2017.
- 506 [65] S. Kaur, J. Cohen, and Z. C. Lipton. Are perceptually-aligned gradients a general property of robust
507 classifiers? In *NeurIPS Workshop: Science Meets Engineering of Deep Learning*, 2019.
- 508 [66] J. Kim and X. Wang. Sensible adversarial learning. *OpenReview*, 2019.
- 509 [67] K. Kireev, M. Andriushchenko, and N. Flammarion. On the effectiveness of adversarial training against
510 common corruptions. *arXiv*, 2021.
- 511 [68] P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Yasunaga,
512 R. L. Phillips, I. Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. *arXiv preprint*
513 *arXiv:2012.07421*, 2020.
- 514 [69] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Technical Report*,
515 2009.
- 516 [70] S. Kundu, M. Nazemi, P. A. Beerel, and M. Pedram. A tunable robust pruning framework through
517 dynamic network rewiring of dnns. *ASP-DAC*, 2021.
- 518 [71] A. Kurakin, I. Goodfellow, S. Bengio, Y. Dong, F. Liao, M. Liang, T. Pang, J. Zhu, X. Hu, C. Xie, et al.
519 Adversarial attacks and defences competition. In *NeurIPS Competition Track*, 2018.
- 520 [72] C. Laidlaw and S. Feizi. Functional adversarial attacks. In *NeurIPS*, 2019.
- 521 [73] C. Laidlaw, S. Singla, and S. Feizi. Perceptual adversarial robustness: Defense against unseen threat
522 models. *arXiv preprint arXiv:2006.12655*, 2020.
- 523 [74] M. Lecuyer, V. Atlidakis, R. Geambasu, D. Hsu, and S. Jana. Certified robustness to adversarial examples
524 with differential privacy. In *2019 IEEE S&P*, 2019.
- 525 [75] L. Li, M. Weber, X. Xu, L. Rimanic, B. Kailkhura, T. Xie, C. Zhang, and B. Li. Tss: Transformation-
526 specific smoothing for robustness certification. In *ACM CCS*, 2021.
- 527 [76] Y. Li, J. Bradshaw, and Y. Sharma. Are generative classifiers more robust to adversarial attacks? In *ICML*,
528 2019.
- 529 [77] X. Ling, S. Ji, J. Zou, J. Wang, C. Wu, B. Li, and T. Wang. Deepsec: A uniform platform for security
530 analysis of deep learning model. In *IEEE S&P*, 2019.

- 531 [78] D. C. Liu and J. Nocedal. On the limited memory bfgs method for large scale optimization. *Mathematical*
532 *programming*, 45(1):503–528, 1989.
- 533 [79] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to
534 adversarial attacks. In *ICLR*, 2018.
- 535 [80] T. Maho, B. Bonnet, T. Furon, and E. L. Merrer. Robic: A benchmark suite for assessing classifiers
536 robustness. *arXiv preprint arXiv:2102.05368*, 2021.
- 537 [81] P. Maini, E. Wong, and J. Z. Kolter. Adversarial robustness against the union of multiple perturbation
538 models. In *ICML*, 2020.
- 539 [82] C. Mao, Z. Zhong, J. Yang, C. Vondrick, and B. Ray. Metric learning for adversarial robustness. *NeurIPS*,
540 2019.
- 541 [83] M. Melis, A. Demontis, M. Pintor, A. Sotgiu, and B. Biggio. secml: A python library for secure and
542 explainable machine learning. *arXiv preprint arXiv:1912.10013*, 2019.
- 543 [84] A. Modas, S.-M. Moosavi-Dezfooli, and P. Frossard. Sparsefool: a few pixels make a big difference. In
544 *CVPR*, 2019.
- 545 [85] S.-M. Moosavi-Dezfooli, A. Fawzi, J. Uesato, and P. Frossard. Robustness via curvature regularization,
546 and vice versa. *CVPR*, 2019.
- 547 [86] M. Mosbach, M. Andriushchenko, T. Trost, M. Hein, and D. Klakow. Logit pairing methods can fool
548 gradient-based attacks. In *NeurIPS 2018 Workshop on Security in Machine Learning*, 2018.
- 549 [87] A. Mustafa, S. Khan, M. Hayat, R. Goecke, J. Shen, and L. Shao. Adversarial defense by restricting the
550 hidden space of deep neural networks. *ICCV*, 2019.
- 551 [88] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with
552 unsupervised feature learning. *Technical Report*, 2011.
- 553 [89] M.-I. Nicolae, M. Sinn, M. N. Tran, B. Buesser, A. Rawat, M. Wistuba, V. Zantedeschi, N. Baracaldo,
554 B. Chen, H. Ludwig, I. Molloy, and B. Edwards. Adversarial robustness toolbox v1.2.0. *arXiv preprint*
555 *arXiv:1807.01069*, 2018.
- 556 [90] T. Pang, K. Xu, Y. Dong, C. Du, N. Chen, and J. Zhu. Rethinking softmax cross-entropy loss for
557 adversarial robustness. *ICLR*, 2020.
- 558 [91] T. Pang, K. Xu, and J. Zhu. Mixup inference: Better exploiting mixup to defend adversarial attacks. In
559 *ICLR*, 2020.
- 560 [92] T. Pang, X. Yang, Y. Dong, K. Xu, H. Su, and J. Zhu. Boosting adversarial training with hypersphere
561 embedding. *NeurIPS*, 2020.
- 562 [93] T. Pang, X. Yang, Y. Dong, H. Su, and J. Zhu. Bag of tricks for adversarial training. *ICLR*, 2021.
- 563 [94] T. Pang, H. Zhang, D. He, Y. Dong, H. Su, W. Chen, J. Zhu, and T.-Y. Liu. Adversarial training with
564 rectified rejection. *arXiv preprint arXiv:2105.14785*, 2021.
- 565 [95] N. Papernot, F. Faghri, N. Carlini, I. Goodfellow, R. Feinman, A. Kurakin, C. Xie, Y. Sharma, T. Brown,
566 A. Roy, A. Matyasko, V. Behzadan, K. Hambardzumyan, Z. Zhang, Y.-L. Juang, Z. Li, R. Sheatsley,
567 A. Garg, J. Uesato, W. Gierke, Y. Dong, D. Berthelot, P. Hendricks, J. Rauber, and R. Long. Technical
568 report on the clevehans v2.1.0 adversarial examples library. *arXiv preprint arXiv:1610.00768*, 2018.
- 569 [96] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and
570 A. Lerer. Automatic differentiation in pytorch. *Technical Report*, 2017.
- 571 [97] C. Qin, J. Martens, S. Goyal, D. Krishnan, K. Dvijotham, A. Fawzi, S. De, R. Stanforth, and P. Kohli.
572 Adversarial robustness through local linearization. *NeurIPS*, 2019.
- 573 [98] A. Raghunathan, S. M. Xie, F. Yang, J. Duchi, and P. Liang. Understanding and mitigating the tradeoff
574 between robustness and accuracy. In *ICML*, 2020.
- 575 [99] M. S. Rahman, M. Imani, N. Mathews, and M. Wright. Mockingbird: Defending against deep-learning-
576 based website fingerprinting attacks with adversarial traces. *IEEE Transactions on Information Forensics*
577 *and Security*, 16:1594–1609, 2020.
- 578 [100] J. Rauber, W. Brendel, and M. Bethge. Foolbox: A python toolbox to benchmark the robustness of
579 machine learning models. In *ICML Reliable Machine Learning in the Wild Workshop*, 2017.
- 580 [101] S.-A. Rebuffi, S. Goyal, D. A. Calian, F. Stimberg, O. Wiles, and T. Mann. Fixing data augmentation to
581 improve adversarial robustness. *arXiv preprint arXiv:2103.01946*, 2021. URL [https://arxiv.org/
582 pdf/2103.01946](https://arxiv.org/pdf/2103.01946).
- 583 [102] B. Recht, R. Roelofs, L. Schmidt, and V. Shankar. Do imagenet classifiers generalize to imagenet? In
584 *ICML*, 2019.
- 585 [103] L. Rice, E. Wong, and J. Z. Kolter. Overfitting in adversarially robust deep learning. In *ICML*, 2020.

- 586 [104] J. Rony, L. G. Hafemann, L. S. Oliveira, I. B. Ayed, R. Sabourin, and E. Granger. Decoupling direction
587 and norm for efficient gradient-based l2 adversarial attacks and defenses. *CVPR*, 2019.
- 588 [105] J. Rony, E. Granger, M. Pedersoli, and I. Ben Ayed. Augmented lagrangian adversarial attacks. *arXiv*
589 *preprint arXiv:2011.11857*, 2020.
- 590 [106] P. Saadatpanah, A. Shafahi, and T. Goldstein. Adversarial attacks on copyright detection systems. In
591 *ICML*, 2020.
- 592 [107] H. Salman, A. Ilyas, L. Engstrom, A. Kapoor, and A. Madry. Do adversarially robust imagenet models
593 transfer better? *NeurIPS*, 2020.
- 594 [108] P. Samangouei, M. Kabkab, and R. Chellappa. Defense-GAN: Protecting classifiers against adversarial
595 attacks using generative models. In *ICLR*, 2018.
- 596 [109] L. Schwinn, R. Raab, A. Nguyen, D. Zanca, and B. Eskofier. Exploring robust misclassifications of neural
597 networks to enhance adversarial attacks. *arXiv preprint arXiv:2105.10304*, 2021.
- 598 [110] V. Sehwag, A. N. Bhagoji, L. Song, C. Sitawarin, D. Cullina, M. Chiang, and P. Mittal. Analyzing
599 the robustness of open-world machine learning. In *12th ACM Workshop on Artificial Intelligence and*
600 *Security*, 2019.
- 601 [111] V. Sehwag, S. Wang, P. Mittal, and S. Jana. Hydra: Pruning adversarially robust neural networks. *NeurIPS*,
602 2020.
- 603 [112] V. Sehwag, S. Wang, P. Mittal, and S. Jana. On pruning adversarially robust neural networks. *NeurIPS*,
604 2020.
- 605 [113] V. Sehwag, S. Mahloujifar, T. Handina, S. Dai, C. Xiang, M. Chiang, and P. Mittal. Improving adversarial
606 robustness using proxy distributions. *arXiv*, 2021.
- 607 [114] A. Shafahi, M. Najibi, A. Ghiasi, Z. Xu, J. Dickerson, C. Studer, L. S. Davis, G. Taylor, and T. Goldstein.
608 Adversarial training for free! *NeurIPS*, 2019.
- 609 [115] S. Shan, E. Wenger, J. Zhang, H. Li, H. Zheng, and B. Y. Zhao. Fawkes: Protecting privacy against
610 unauthorized deep learning models. In *29th USENIX Security Symposium (USENIX Security 20)*, pages
611 1589–1604, 2020.
- 612 [116] C. Shi, C. Holtz, and G. Mishne. Online adversarial purification based on self-supervised learning.
613 In *International Conference on Learning Representations*, 2021. URL [https://openreview.net/
614 forum?id=_i3ASp12WS](https://openreview.net/forum?id=_i3ASp12WS).
- 615 [117] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine
616 learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18. IEEE, 2017.
- 617 [118] M. Singh, A. Sinha, N. Kumari, H. Machiraju, B. Krishnamurthy, and V. N. Balasubramanian. Harnessing
618 the vulnerability of latent layers in adversarially trained models. *IJCAI*, 2019.
- 619 [119] C. Sitawarin, S. Chakraborty, and D. Wagner. Improving adversarial robustness through progressive
620 hardening. *arXiv*, 2020.
- 621 [120] L. Song and P. Mittal. Systematic evaluation of privacy risks of machine learning models. In *30th USENIX*
622 *Security Symposium (USENIX Security 21)*, 2021.
- 623 [121] L. Song, R. Shokri, and P. Mittal. Privacy risks of securing machine learning models against adversarial
624 examples. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications*
625 *Security*, pages 241–257, 2019.
- 626 [122] L. Song, V. Sehwag, A. N. Bhagoji, and P. Mittal. A critical evaluation of open-world machine learning.
627 *arXiv preprint arXiv:2007.04391*, 2020.
- 628 [123] C. Szegedy, W. Zaremba, I. Sutskever, D. E. Joan Bruna, I. Goodfellow, and R. Fergus. Intriguing
629 properties of neural networks. In *ICLR*, 2013.
- 630 [124] L. Tao, L. Feng, J. Yi, S.-J. Huang, and S. Chen. Provable defense against delusive poisoning. *arXiv*
631 *preprint arXiv:2102.04716*, 2021.
- 632 [125] R. Taori, A. Dave, V. Shankar, N. Carlini, B. Recht, and L. Schmidt. Measuring robustness to natural
633 distribution shifts in image classification. *arXiv preprint arXiv:2007.00644*, 2020.
- 634 [126] V. Tjeng, K. Xiao, and R. Tedrake. Evaluating robustness of neural networks with mixed integer
635 programming. In *ICLR*, 2019.
- 636 [127] F. Tramèr and D. Boneh. Adversarial training and robustness for multiple perturbations. In *NeurIPS*,
637 2019.
- 638 [128] F. Tramèr, P. Dupré, G. Rusak, G. Pellegrino, and D. Boneh. Adversarial: Perceptual ad blocking meets
639 adversarial machine learning. In *ACM SIGSAC CCS*, 2019.

- 640 [129] F. Tramèr, N. Carlini, W. Brendel, and A. Madry. On adaptive attacks to adversarial example defenses. In
641 *NeurIPS*, 2020.
- 642 [130] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness may be at odds with accuracy.
643 In *ICLR*, 2019.
- 644 [131] J. Uesato, J.-B. Alayrac, P.-S. Huang, R. Stanforth, A. Fawzi, and P. Kohli. Are labels required for
645 improving adversarial robustness? *NeurIPS*, 2019.
- 646 [132] F. Utrera, E. Kravitz, N. B. Erichson, R. Khanna, and M. W. Mahoney. Adversarially-trained deep nets
647 transfer better. *arXiv preprint arXiv:2007.05869*, 2020.
- 648 [133] D. Wang, A. Ju, E. Shelhamer, D. Wagner, and T. Darrell. Fighting gradients with gradients: Dynamic
649 defenses against adversarial attacks. *arXiv preprint arXiv:2105.08714*, 2021.
- 650 [134] J. Wang and H. Zhang. Bilateral adversarial training: Towards fast training of more robust models against
651 adversarial attacks. *ICCV*, 2019.
- 652 [135] Y. Wang, D. Zou, J. Yi, J. Bailey, X. Ma, and Q. Gu. Improving adversarial robustness requires revisiting
653 misclassified examples. *ICLR*, 2020.
- 654 [136] T.-W. Weng, H. Zhang, H. Chen, Z. Song, C.-J. Hsieh, D. Boning, I. S. Dhillon, and L. Daniel. Towards
655 fast computation of certified robustness for relu networks. In *ICML*, 2018.
- 656 [137] E. Wong and J. Z. Kolter. Learning perturbation sets for robust machine learning. *arXiv preprint*
657 *arXiv:2007.08450*, 2020.
- 658 [138] E. Wong and Z. Kolter. Provable defenses against adversarial examples via the convex outer adversarial
659 polytope. *ICML*, 2018.
- 660 [139] E. Wong, F. R. Schmidt, and J. Z. Kolter. Wasserstein adversarial examples via projected sinkhorn
661 iterations. In *ICML*, 2019.
- 662 [140] E. Wong, L. Rice, and J. Z. Kolter. Fast is better than free: Revisiting adversarial training. *ICLR*, 2020.
- 663 [141] B. Wu, J. Chen, D. Cai, X. He, and Q. Gu. Do wider neural networks really help adversarial robustness?
664 *arXiv*, 2020.
- 665 [142] D. Wu, S. tao Xia, and Y. Wang. Adversarial weight perturbation helps robust generalization. *NeurIPS*,
666 2020.
- 667 [143] C. Xiao, P. Zhong, and C. Zheng. Enhancing adversarial defense by k-winners-take-all. *ICLR*, 2020.
- 668 [144] C. Xie, M. Tan, B. Gong, J. Wang, A. L. Yuille, and Q. V. Le. Adversarial examples improve image
669 recognition. In *CVPR*, 2020.
- 670 [145] C. Xu, X. Li, and M. Yang. An orthogonal classifier for improving the adversarial robustness of neural
671 networks. *arXiv preprint arXiv:2105.09109*, 2021.
- 672 [146] H. Xu, X. Liu, Y. Li, and J. Tang. To be robust or to be fair: Towards fairness in adversarial training.
673 *arXiv preprint arXiv:2010.06121*, 2020.
- 674 [147] Y. Yang, G. Zhang, D. Katabi, and Z. Xu. Me-net: Towards effective adversarial robustness with matrix
675 estimation. In *ICML*, 2019.
- 676 [148] Y.-Y. Yang, C. Rashtchian, H. Zhang, R. Salakhutdinov, and K. Chaudhuri. A closer look at accuracy vs.
677 robustness. *Advances in Neural Information Processing Systems*, 33, 2020.
- 678 [149] Y. Yu, Z. Yang, E. Dobriban, J. Steinhardt, and Y. Ma. Understanding generalization in adversarial
679 training via the bias-variance decomposition. *arXiv preprint arXiv:2103.09947*, 2021.
- 680 [150] S. Zagoruyko and N. Komodakis. Wide residual networks. In *BMVC*, 2016.
- 681 [151] D. Zhang, T. Zhang, Y. Lu, Z. Zhu, and B. Dong. You only propagate once: Accelerating adversarial
682 training via maximal principle. *NeurIPS*, 2019.
- 683 [152] H. Zhang and J. Wang. Defense against adversarial attacks using feature scattering-based adversarial
684 training. *NeurIPS*, 2019.
- 685 [153] H. Zhang and W. Xu. Adversarial interpolation training: A simple approach for improving model
686 robustness. *OpenReview*, 2019.
- 687 [154] H. Zhang, Y. Yu, J. Jiao, E. P. Xing, L. E. Ghaoui, and M. I. Jordan. Theoretically principled trade-off
688 between robustness and accuracy. In *ICML*, 2019.
- 689 [155] J. Zhang, X. Xu, B. Han, G. Niu, L. Cui, M. Sugiyama, and M. Kankanhalli. Attacks which do not kill
690 training make adversarial learning stronger. *ICML*, 2020.
- 691 [156] J. Zhang, J. Zhu, G. Niu, B. Han, M. Sugiyama, and M. Kankanhalli. Geometry-aware instance-reweighted
692 adversarial training. *ICLR*, 2021.

- 693 [157] J. Zhong, X. Liu, and C.-J. Hsieh. Improving the speed and quality of gan by adversarial training. *arXiv preprint arXiv:2008.03364*, 2020.
694
695 [158] C. Zhu, Y. Cheng, Z. Gan, S. Sun, T. Goldstein, and J. Liu. Freelib: Enhanced adversarial training for
696 natural language understanding. In *ICLR*, 2019.

697 Checklist

- 698 1. For all authors...
- 699 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
700 contributions and scope? [Yes] We clearly describe all aspects of the benchmark in
701 Sec. 3, Sec. 4 and in the Appendix.
- 702 (b) Did you describe the limitations of your work? [Yes] We discuss that we only perform
703 a standardized evaluation and do not evaluate models against adaptive attacks, although
704 we accept third-party evaluations based on adaptive attacks as mentioned in Sec. 3.1.
705 Also, our current set of leaderboards can be seen as a limitation, so in Sec. 5 we describe
706 how we plan to expand the benchmark to new threat models and datasets.
- 707 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 708 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
709 them? [Yes]
- 710 2. If you are including theoretical results...
- 711 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 712 (b) Did you include complete proofs of all theoretical results? [N/A]
- 713 3. If you ran experiments (e.g. for benchmarks)...
- 714 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
715 mental results (either in the supplemental material or as a URL)? [Yes] See Sec. 3.2
716 and the Model Zoo page.
- 717 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
718 were chosen)? [N/A]
- 719 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
720 ments multiple times)? [Yes] See Sec. C.
- 721 (d) Did you include the total amount of compute and the type of resources used (e.g., type
722 of GPUs, internal cluster, or cloud provider)? [Yes] We give in Sec. C an example of
723 the computational time and infrastructure used to run AutoAttack on several models.
724 However, we could not provide such details for all models in the benchmark since
725 those were collected over time with different resources.
- 726 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 727 (a) If your work uses existing assets, did you cite the creators? [Yes] See Sec. E, and in
728 the text we cite the authors of the datasets and algorithms we use.
- 729 (b) Did you mention the license of the assets? [Yes] See Sec. A.
- 730 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
731 We provide code and data in the Model Zoo page.
- 732 (d) Did you discuss whether and how consent was obtained from people whose data you’re
733 using/curating? [Yes] See Sec. 3.2 and Sec. A.
- 734 (e) Did you discuss whether the data you are using/curating contains personally identifiable
735 information or offensive content? [N/A]
- 736 5. If you used crowdsourcing or conducted research with human subjects...
- 737 (a) Did you include the full text of instructions given to participants and screenshots, if
738 applicable? [N/A]
- 739 (b) Did you describe any potential participant risks, with links to Institutional Review
740 Board (IRB) approvals, if applicable? [N/A]
- 741 (c) Did you include the estimated hourly wage paid to participants and the total amount
742 spent on participant compensation? [N/A]