

A Licenses

The code used for benchmarking is released under MIT license. The code of AutoAttack [26] that our benchmark relies on has been released under the MIT license as well. The classifiers in the Model Zoo are added according to the permission given by the authors with the license they choose: most of the models have MIT license, other have more restrictive ones such as Attribution-NonCommercial-ShareAlike 4.0 International, Apache License 2.0, BSD 3-Clause License. The details can be found at <https://github.com/RobustBench/robustbench/blob/master/LICENSE>. The CIFAR-10 and CIFAR-100 datasets [69] are obtained via the PyTorch loaders [96], while CIFAR-10-C and CIFAR-100-C [53], with the common corruptions, are downloaded from the official release (see <https://zenodo.org/record/2535967#.YLYf9agzaUk> and <https://zenodo.org/record/3555552#.YLYeJagzaUk>).

B Maintenance plan

Here we discuss the main aspects of maintaining RobustBench and the costs associated with it:

- **Hosting the website** (<https://robustbench.github.io/>): we host our leaderboard using GitHub pages² which is a free service.
- **Hosting the library** (<https://github.com/RobustBench/robustbench>): the code of our library is hosted on GitHub³ which offers the basic features that we need to maintain the library for free.
- **Hosting the models**: to ensure the availability of the models from the Model Zoo, we host them in our own cloud storage on Google Drive⁴. At the moment, they take around 24 GB of space which fits into the 100 GB storage plan that costs 2 USD per month.
- **Running evaluations**: we run all evaluations on the GPU servers that are available to our research groups which incurs no extra costs.

Moreover, as we mention in the outlook (Sec. 5), we also plan to expand the benchmark to new datasets and threat models which can slightly increase the required maintenance costs since we may need to upgrade the storage plan. We also expect the benchmark to be community-driven and to encourage this we have provided instructions⁵ on how to submit new entries to the leaderboard and to the Model Zoo.

C Reproducibility and runtime

Here we discuss the main aspects of the reproducibility of the benchmark.

First of all, the code to run the benchmark on a given model is available in our repository, and an example of how to run it is given in the README file. The installation instructions are also provided in the README file and the requirements will be installed automatically. We only leave to the user the installation of the PyTorch and torchvision libraries to allow for the installation of the most appropriate versions for the user's hardware (e.g. CUDA vs CPU versions). To satisfy other points from the reproducibility checklist⁶ which are applicable to our benchmark, we also discuss next the variability of the robust accuracy over random seeds and the average runtime of the benchmark.

Evaluation of the accuracy on common corruptions [53] is deterministic if we do not take into account non-deterministic operations on computational accelerators such as GPUs⁷ which, however, do not affect the resulting accuracy. On the other hand, robustness evaluation using AutoAttack has an element of randomness since it relies on random initializations of the starting points and also on the randomness in the update of the Square Attack [4]. To show the effect of randomness on the

²<https://pages.github.com/>

³<http://github.com/>

⁴<https://www.google.com/drive/>

⁵<https://github.com/RobustBench/robustbench#adding-a-new-model>

⁶<https://www.cs.mcgill.ca/~jpineau/ReproducibilityChecklist.pdf>

⁷<https://pytorch.org/docs/stable/notes/randomness.html>

robust accuracy given by AutoAttack, we repeat evaluation over four random seeds on three models available in the Model Zoo from different threat models. In Table 2 we report the average robust accuracy with its standard deviation. We observe that different seeds lead to very similar results. Moreover, we indicate the runtime of each evaluation, which is largely influenced by the size of the model and the computing infrastructure (every run uses a single Tesla V100 GPU, and we set a batch size of 500). Moreover, less robust models require less time for evaluation which is due to the fact that AutoAttack does not further attack a point if an adversarial example is already found by some preceding attack in the ensemble.

Table 2: Statistics about the standardized evaluation with AutoAttack when repeated for four random seeds. We can see that the robust accuracy has very small fluctuations. We also report the runtime for the different models which is much smaller for less robust models.

Dataset	Leaderboard	Paper	Architecture	Clean acc.	Robust acc.	Time
CIFAR-10	ℓ_∞	Gowal et al. [46]	WRN-28-10	89.48%	$62.82\% \pm 0.016$	11.8 h
CIFAR-10	ℓ_2	Rebuffi et al. [100]	WRN-28-10	91.79%	$78.80\% \pm 0.000$	15.1 h
CIFAR-100	ℓ_∞	Wu et al. [140]	WRN-34-10	60.38%	$28.84\% \pm 0.018$	6.6 h

D Additional analysis

In this section, we show more results on different datasets and/or threat models and discuss some implementation details related to the analysis from Sec. 4. We also additionally analyze the *smoothness* and *transferability* properties of the models from the Model Zoo.

Progress on adversarial defenses. As done in the main part for the ℓ_∞ -robust models on CIFAR-10, we show here the same statistics but for ℓ_2 -robust models on CIFAR-10 in Fig. 8 and for ℓ_∞ -robust models on CIFAR-100 in Fig. 9. We observe a few differences compared to the ℓ_∞ -robust models on CIFAR-10 reported in Fig. 2. First of all, the amount of robustness overestimation is not large and in particular there are no models that have *zero* robust accuracy. Second, we can see that the best ℓ_2 -robust models on CIFAR-10 has even higher standard accuracy than a standard model (95.74% vs 94.78%) while having a significantly higher robust accuracy (82.32% vs 0.00%) and leaving a relatively small gap between the standard and robust accuracy. Finally, we note that the progress on the ℓ_∞ -threat model on CIFAR-100 is more recent and there are only a few published papers that report adversarial robustness on this dataset.

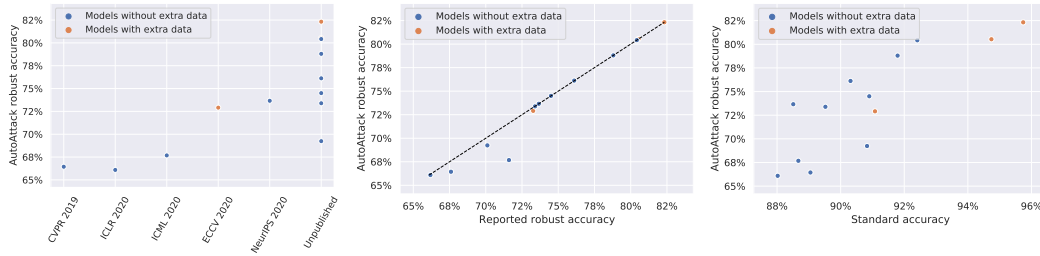


Figure 8: Visualization of the robustness and accuracy of 13 CIFAR-10 models from the RobustBench ℓ_2 -leaderboard. Robustness is evaluated using ℓ_2 -perturbations with $\varepsilon_2 = 0.5$.

Robustness across distribution shifts. We measure robust accuracy on various distribution shifts using four dataset, namely CIFAR-10, CINIC-10, CIFAR-10.1, and CIFAR-10-C. In particular, we compute the robust accuracy in the same threat model as for the original CIFAR-10 dataset, and report the results in Fig. 10. Interestingly, one can observe that ℓ_p adversarial robustness is maintained under the distribution shifts, and it highly correlates with the robustness on the dataset the models were trained on (i.e. CIFAR-10).

Calibration. We compute the expected calibration error (ECE) using the code of [47]. We use their default settings to compute the calibration error which includes, in particular, binning of the probability range onto 15 equally-sized bins. However, we use our own implementation of the



Figure 9: Visualization of the robustness and accuracy of 12 CIFAR-100 models from the RobustBench ℓ_∞ -leaderboard. Robustness is evaluated using ℓ_∞ -perturbations with $\varepsilon_\infty = 8/255$.

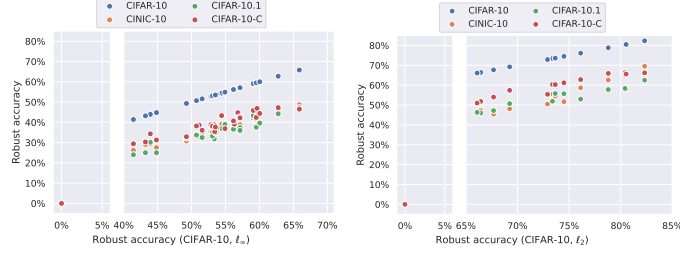


Figure 10: Robust accuracy of the robust classifiers, trained against ℓ_∞ and ℓ_2 threat model, respectively, from our Model Zoo on various distribution shifts. The data points with 0% robust accuracy correspond to a standardly trained model.

799 temperature rescaling algorithm which is close to that of [7]. Since optimization of the ECE over the
800 softmax temperature is a simple *one-dimensional* optimization problem, we can solve it efficiently
801 using a grid search. Moreover, the advantage of performing a grid search is that we can optimize
802 directly the metric of interest, i.e. ECE, instead of the cross-entropy loss as in [47] who had to
803 rely on a differentiable loss since they used LBFGS [78] to optimize the temperature. We perform
804 a grid search over the interval $t \in [0.001, 1.0]$ with a grid step 0.001 and we test both t and $1/t$
805 temperatures. Moreover, we check that for all models the optimal temperature t is situated not at the
806 boundary of the grid.

807 We show additional calibration results for ℓ_2 -robust models in Fig. 11. The overall trend of the ECE
808 is the same as for ℓ_∞ -robust models: most of the ℓ_2 models are underconfident (since the optimal
809 temperature is less than one) and lead to worse calibration before and after temperature rescaling.
810 The main difference compared to the ℓ_∞ threat model is that among the ℓ_2 models there are two
811 models that are *better-calibrated*: one before (Engstrom et al. [35] with 1.41% ECE vs 3.71% ECE
812 of the standard model) and one after (Gowal et al. [46] with 1.00% ECE vs 1.11% ECE of the
813 standard model) temperature rescaling. Moreover, we can see that similarly to the ℓ_∞ case, the only
814 overconfident models are either the standard one or models that maximize the margin instead of using
norm-bounded perturbations, i.e. Ding et al. [32] and Rony et al. [103].

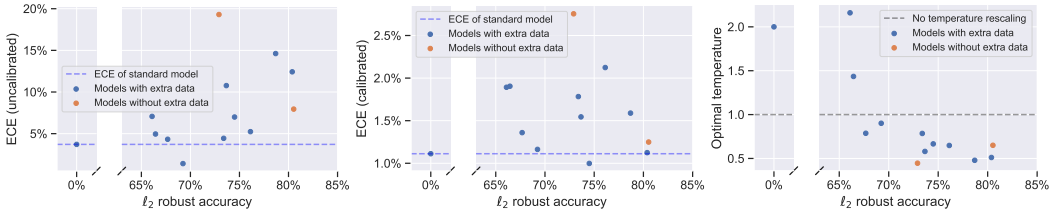


Figure 11: Expected calibration error (ECE) before (**left**) and after (**middle**) temperature rescaling, and the optimal rescaling temperature (**right**) for the ℓ_2 -robust models.

815

816 **Out-of-distribution detection.** Fig. 12 complements Fig. 5 and shows the ability of ℓ_2 -robust models
817 trained on CIFAR-10 to distinguish inputs from other datasets (CIFAR-100, SVHN, Describable
818 Textures). We find that ℓ_2 robust models have in general comparable OOD detection performance

819 to standardly trained models, while the model by Augustin et al. [7] achieves even better perfor-
820 mance since their approach explicitly optimizes both robust accuracy and worst-case OOD detection
performance.

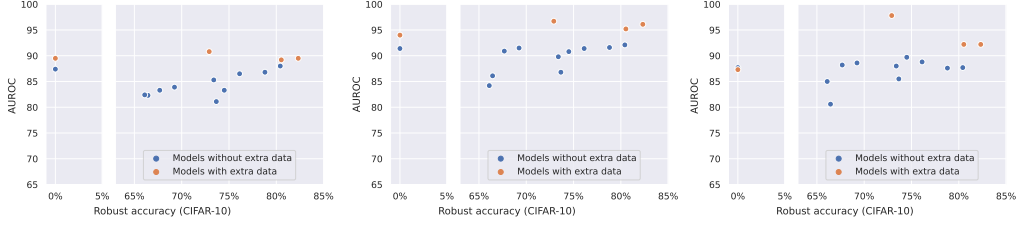


Figure 12: Visualization of the quality of OOD detection (higher AUROC is better) for the ℓ_2 -robust models on three different OOD datasets: CIFAR-100 (left), SVHN (middle), Describable Textures (right).

821

822 **Fairness in robustness.** We report the results about fairness for robust models in the ℓ_2 -threat model
823 in Fig. 13, similarly to what done for ℓ_∞ above. We see that the difference in robustness among
824 classes is similar to what observed for the ℓ_∞ models. Also, the RSD of robustness over classes
825 decreases, which indicates that the disparity among subgroups is reduced, as the average robust
826 accuracy improves. To compute the robustness for the experiments about fairness we used APGD on
827 the targeted DLR loss [26] with 3 target classes and 20 iterations each on the whole test set. Note
828 that even with this smaller budget we achieve results very close to that of the full evaluation, with an
average difference smaller than 0.5%.

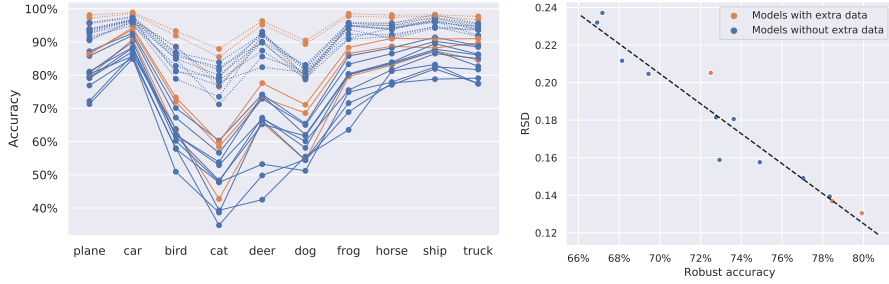


Figure 13: **Left:** classwise standard (dotted lines) and robust (solid) accuracy of ℓ_2 -robust models. **Right:** relative standard deviation (RSD) of robust accuracy over classes vs its average.

829

830 **Privacy leakage.** We use membership inference accuracy, referred to as inference accuracy, as a
831 measure of the leakage of training data details from pre-trained neural networks. It measures how
832 successfully we can identify whether a particular sample was present in the training set. We closely
833 follow the methodology described in Song and Mittal [118] to calculate inference accuracy. In
834 particular, we measure the confidence in the correct class for each input image with a pre-trained
835 classifier. We measure the confidence for both training and test set images and calculate the maximum
836 classification accuracy between train and test images based on the confidence values. We refer to this
837 accuracy as *inference accuracy using confidence*. We also follow the recommendation from Song
838 et al. [119] where they show that adversarial examples are more successful in estimating inference
839 accuracy on robust networks. In our experiments, we also find that using adversarial examples leads
840 to higher inference accuracy than benign images (Figure 14). We also find that robust networks in
841 the ℓ_2 threat model have relatively higher inference accuracy than robust networks in the ℓ_∞ threat
842 model.

843 A key reason behind privacy leakage through membership inference is that deep neural networks often
844 end up overfitting on the training data. One standard metric to measure overfitting is the generalization
845 gap between train and test set. Naturally, this difference in the accuracy on the train and test set is
846 the baseline of inference accuracy. We refer to it as *inference accuracy using label* and report it in
847 Figure 15. We consider both benign and adversarial images. When using benign images, we find
848 confidence information does lead to higher inference accuracy than using only labels. However, with

adversarial examples, which achieve higher inference accuracy than benign images, we find that inference accuracy based on confidence information closely follows the inference accuracy calculate from labels.

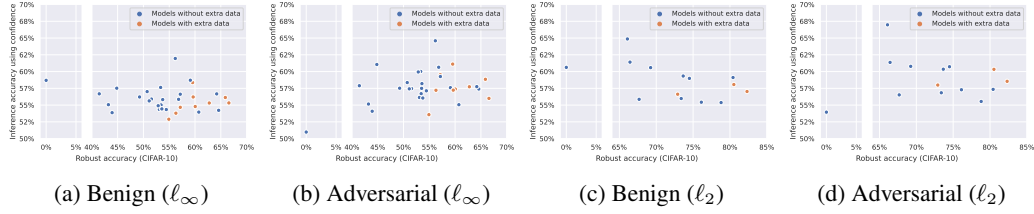


Figure 14: **Comparing privacy leakage of different networks.** We compare membership inference accuracy from benign and adversarial images across both ℓ_∞ and ℓ_2 threat model.

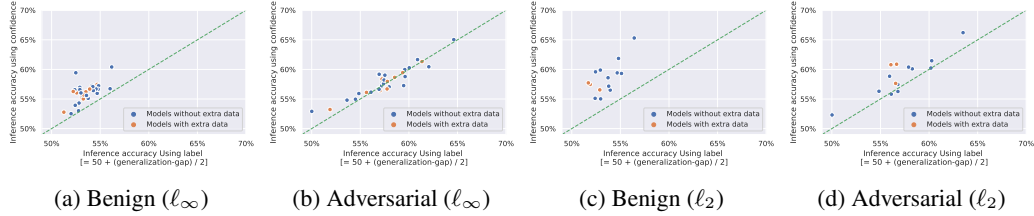


Figure 15: **Comparing privacy leakage with different output statistics.** We measure privacy leakage using membership inference accuracy, i.e., classification success between train and test set. We measure it using two baselines 1) based on correct prediction i.e., using predicted class label and 2) based on classification confidence in correct class. We also measure it using both benign and adversarial images.

Smoothness. Previous work [146] has shown that smoothness of a model, together with enough separation between the classes of the dataset for which it is trained, is necessary to achieve both natural and robust accuracy. They use local Lipschitzness as a measure for model smoothness, and observe empirically that robust models are more smoother than models trained in a standard way. Our Model Zoo enables us to check this fact empirically on a wider range of robust models, trained with a more diverse set of techniques, in particular with and without extra training data. Moreover, as we have access to the model internals, we can also compute local Lipschitzness of the model up to arbitrary layers, to see how smoothness changes between layers.

We compute local Lipschitzness using projected gradient descent (PGD) on the following optimization problem:

$$L = \frac{1}{N} \sum_{i=1}^N \max_{\substack{x_1: \|x_1 - x_i\|_\infty \leq \varepsilon, \\ x_2: \|x_2 - x_i\|_\infty \leq \varepsilon}} \frac{\|f(x_1) - f(x_2)\|_1}{\|x_1 - x_2\|_\infty}, \quad (2)$$

where x_i represents each sample around which we compute local Lipschitzness, N is the number of samples across which we average ($N = 256$ in all our experiments), and f represents the function whose Lipschitz constant we compute. As mentioned above, this function can be either the full model, or the model up to an arbitrary intermediate layer.

Since the models can have similar smoothness behavior, but at a different scale, we also consider normalizing the models outputs. One such normalization we use is given by the projection of the model outputs on the unit ℓ_2 ball. This normalization aims at capturing the angular change of the output, instead of taking in consideration also its magnitude. We compute the “angular” version of the Lipschitz constant as

$$L = \frac{1}{N} \sum_{i=1}^N \max_{\substack{x_1: \|x_1 - x_i\|_\infty \leq \varepsilon, \\ x_2: \|x_2 - x_i\|_\infty \leq \varepsilon}} \frac{\left\| \frac{f(x_1)}{\|f(x_1)\|_2} - \frac{f(x_2)}{\|f(x_2)\|_2} \right\|_1}{\|x_1 - x_2\|_\infty}. \quad (3)$$

For both variations of Lipschitzness, we compute it with $\varepsilon = 8/255$, running the PDG-like procedure for 50 steps, with a step size of $\varepsilon/5$.



Figure 16: **Lipschitzness**. Computation of the local Lipschitz constant of the WRN-28-10 ℓ_∞ -robust models in our Model Zoo with $\varepsilon = 8/255$. The color coding of the models is the same across both figures. For the correspondence between model IDs (shown in the legend) and papers that introduced them, see Appendix E.

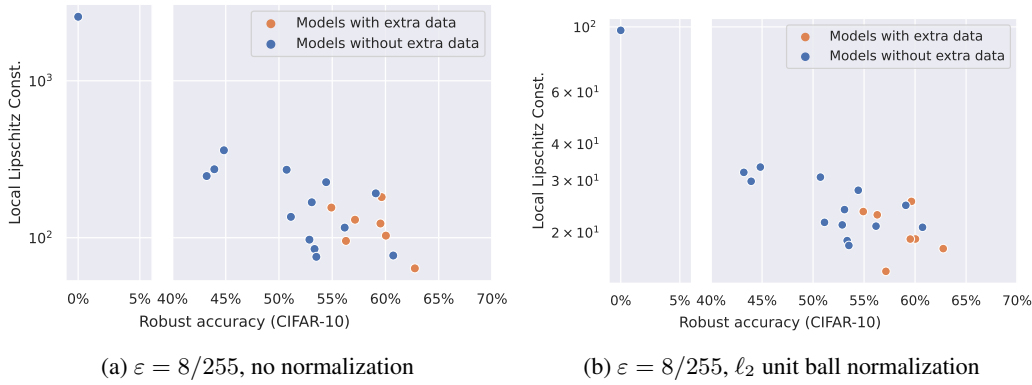


Figure 17: **Lipschitzness vs Robustness**. Local Lipschitz constant of the output layer vs. robust accuracy of a subset of the ℓ_∞ -robust models in our Model Zoo.

In Fig. 16 we compute the layerwise Lipschitzness for all ℓ_∞ models trained on CIFAR-10 from the Model Zoo that have the WRN-28-10 architecture. We observe that the standard model is the least smooth at all the layers, and that all the robustly trained models are smoother. Moreover, we can notice that in Fig. 16a there are two models in the middle ground: these are the models by Gowal et al. [46] and Rebuffi et al. [100], which are the most robust ones, up to the last layer, the smoothest. Nonetheless, in the middle layers, they are the second and third least smooth, according to the unnormalized local Lipschitzness. This can be due to the different activation function used in these models (Swish vs ReLU). For this reason, we also compute “angular” Lipschitzness according to Eq. 3. Indeed, in Fig. 16b, all the robust models are in the same order of magnitude at all layers.

Finally, we also show the Lipschitz constants of the output layer for a larger set of ℓ_∞ models from the Model Zoo that are not restricted to the same architecture. We plot the Lipschitz constant vs. the robust accuracy for these models in Fig. 17. We see that there is a clear relationship between robust accuracy and Lipschitzness, hence confirming the findings of Yang et al. [146].

Transferability. We generate adversarial examples for a network, referred to as source network, and measure robust accuracy of every other network, referred to as target network, from the model zoo on them. We also include additional non-robust models⁸, to name a few, VGG19, ResNet18, and DenseNet121, in our analysis. We consider both ten step PGD attack and FGSM attack to generate adversarial examples as two transferability baselines commonly used in the literature.

We present our results in Figure 18, 19 where the correspondence between model IDs and papers that introduced them can be found in Appendix E. We find that transferability to each robust target network follows a similar trend where adversarial examples transfer equally well from another robust networks. Though slight worse than robust network, adversarial example from non-robust network

⁸We train then for 200 epochs and achieve 93-95% clean accuracy for all networks on the CIFAR-10 dataset.

also transfer equally well to robust networks. We observe a strong transferability among non-robust networks with adversarial examples generated from PGD attacks. Adversarial examples generated using the FGSM attack also transfer successfully. However, they achieve lower robust accuracy on the target network. Intriguingly, we observe the weakest transferability from a robust to a non-robust network. This observation holds for all robust source networks across both FGSM and PGD-attack in both ℓ_∞ and ℓ_2 threat model.

E Leaderboards

We here report the details of all the models included in the various leaderboards, for the ℓ_∞ -, ℓ_2 -threat models and common corruptions. In particular, we show for each model the clean accuracy, robust accuracy (either on adversarial attacks or corrupted images), whether additional data is used for training, the architecture used, the venue at which it appeared and, if available, the identifier in the Model Zoo (which is also used in some of the experiments in Sec. D).

Table 3: Leaderboard for the ℓ_∞ -threat model, CIFAR-10.

Model	Clean	AA	Extra data	Architecture	Venue	Model Zoo ID
1 Rebuffi et al. [100]	92.23	66.56	Y	WRN-70-16	arXiv, Mar 2021	Rebuffi2021Fixing_70_16_cutmix_extra
2 Gowal et al. [46]	91.10	65.87	Y	WRN-70-16	arXiv, Oct 2020	Gowal2020Uncovering_70_16_extra
3 Rebuffi et al. [100]	88.50	64.58	N	WRN-106-16	arXiv, Mar 2021	Rebuffi2021Fixing_106_16_cutmix_ddpm
4 Rebuffi et al. [100]	88.54	64.20	N	WRN-70-16	arXiv, Mar 2021	Rebuffi2021Fixing_70_16_cutmix_ddpm
5 Gowal et al. [46]	89.48	62.76	Y	WRN-28-10	arXiv, Oct 2020	Gowal2020Uncovering_28_10_extra
6 Rebuffi et al. [100]	87.33	60.73	N	WRN-28-10	arXiv, Mar 2021	Rebuffi2021Fixing_28_10_cutmix_ddpm
7 Wu et al. [139]	87.67	60.65	Y	WRN-34-15	arXiv, Oct 2020	N/A
8 Wu et al. [140]	88.25	60.04	Y	WRN-28-10	NeurIPS 2020	Wu2020Adversarial_extra
9 Zhang et al. [154]	89.36	59.64	Y	WRN-28-10	ICLR 2021	Zhang2020Geometry
10 Carmon et al. [18]	89.69	59.53	Y	WRN-28-10	NeurIPS 2019	Carmon2019Unlabeled
11 Sehwal et al. [112]	85.85	59.09	N	WRN-34-10	arXiv, Apr 2021	Sehwal2021Proxy
12 Gowal et al. [46]	85.29	57.14	N	WRN-70-16	arXiv, Oct 2020	Gowal2020Uncovering_70_16
13 Sehwal et al. [110]	88.98	57.14	Y	WRN-28-10	NeurIPS 2020	Sehwal2020Hydra
14 Gowal et al. [46]	85.64	56.82	N	WRN-34-20	arXiv, Oct 2020	Gowal2020Uncovering_34_20
15 Wang et al. [133]	87.50	56.29	Y	WRN-28-10	ICLR 2020	Wang2020Improving
16 Wu et al. [140]	85.36	56.17	N	WRN-34-10	NeurIPS 2020	Wu2020Adversarial
17 Uesato et al. [129]	86.46	56.03	Y	WRN-28-10	NeurIPS 2019	N/A
18 Hendrycks et al. [55]	87.11	54.92	Y	WRN-28-10	ICML 2019	Hendrycks2019Using
19 Sehwal et al. [112]	84.38	54.43	N	ResNet-18	arXiv, Apr 2021	Sehwal2021Proxy_R18
20 Pang et al. [93]	86.43	54.39	N	WRN-34-20	ICLR 2021	N/A
21 Pang et al. [92]	85.14	53.74	N	WRN-34-20	NeurIPS 2020	Pang2020Boosting
22 Cui et al. [28]	88.70	53.57	N	WRN-34-20	arXiv, Nov 2020	Cui2020Learnable_34_20
23 Zhang et al. [153]	84.52	53.51	N	WRN-34-10	ICML 2020	Zhang2020Attacks
24 Rice et al. [102]	85.34	53.42	N	WRN-34-20	ICML 2020	Rice2020Overfitting
25 Huang et al. [58]	83.48	53.34	N	WRN-34-10	NeurIPS 2020	Huang2020Self
26 Zhang et al. [152]	84.92	53.08	N	WRN-34-10	ICML 2019	Zhang2019Theoretically
27 Cui et al. [28]	88.22	52.86	N	WRN-34-10	arXiv, Nov 2020	Cui2020Learnable_34_10
28 Qin et al. [97]	86.28	52.84	N	WRN-40-8	NeurIPS 2019	N/A
29 Chen et al. [21]	86.04	51.56	N	ResNet-50	CVPR 2020	Chen2020Adversarial
30 Chen et al. [20]	85.32	51.12	N	WRN-34-10	arXiv, Oct 2020	Chen2020Efficient
31 Sitawarin et al. [117]	86.84	50.72	N	WRN-34-10	arXiv, Mar 2020	Sitawarin2020Improving
32 Engstrom et al. [35]	87.03	49.25	N	ResNet-50	GitHub, Oct 2019	Engstrom2019Robustness
33 Singh et al. [116]	87.80	49.12	N	WRN-34-10	IJCAI 2019	N/A
34 Mao et al. [82]	86.21	47.41	N	WRN-34-10	NeurIPS 2019	N/A
35 Zhang et al. [149]	87.20	44.83	N	WRN-34-10	NeurIPS 2019	Zhang2019You
36 Madry et al. [79]	87.14	44.04	N	WRN-34-10	ICLR 2018	N/A
37 Andriushchenko et al. [3]	79.84	43.93	N	ResNet-18	NeurIPS 2020	Andriushchenko2020Understanding
38 Pang et al. [90]	80.89	43.48	N	ResNet-32	ICLR 2020	N/A
39 Wong et al. [138]	83.34	43.21	N	ResNet-18	ICLR 2020	Wong2020Fast
40 Shafahi et al. [113]	86.11	41.47	N	WRN-34-10	NeurIPS 2019	N/A
41 Ding et al. [32]	84.36	41.44	N	WRN-28-4	ICLR 2020	Ding2020MMA
42 Kundu et al. [70]	87.32	40.41	N	ResNet-18	ASP-DAC 2021	N/A
43 Atzmon et al. [6]	81.30	40.22	N	ResNet-18	NeurIPS 2019	N/A
44 Moosavi-Dezfooli et al. [85]	83.11	38.50	N	ResNet-18	CVPR 2019	N/A
45 Zhang and Wang [150]	89.98	36.64	N	WRN-28-10	NeurIPS 2019	N/A
46 Zhang and Xu [151]	90.25	36.45	N	WRN-28-10	OpenReview, Sep 2019	N/A
47 Jang et al. [59]	78.91	34.95	N	ResNet-20	ICCV 2019	N/A
48 Kim and Wang [66]	91.51	34.22	N	WRN-34-10	OpenReview, Sep 2019	N/A
49 Wang and Zhang [132]	92.80	29.35	N	WRN-28-10	ICCV 2019	N/A
50 Xiao et al. [141]	79.28	18.50	N	DenseNet-121	ICLR 2020	N/A
51 Jin and Rinard [60]	90.84	1.35	N	ResNet-18	arXiv, Mar 2020	N/A
52 Mustafa et al. [87]	89.16	0.28	N	ResNet-110	ICCV 2019	N/A
53 Chan et al. [19]	93.79	0.26	N	WRN-34-10	ICLR 2020	N/A
54 Alfarra et al. [2]	91.03	0.00	N	WRN-28-10	arXiv, Jun 2020	N/A
55 Standard	94.78	0.0	N	WRN-28-10	N/A	Standard

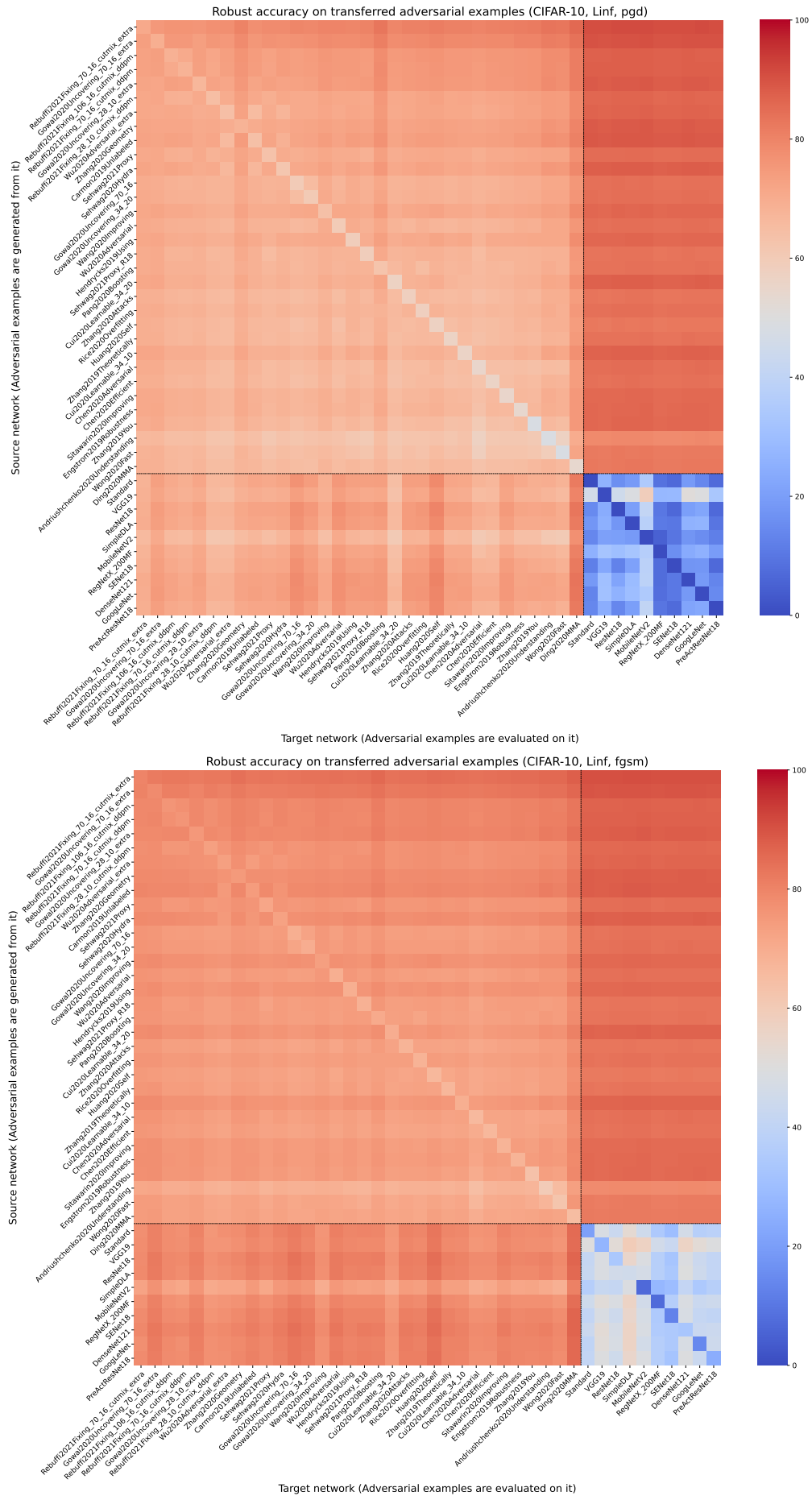


Figure 18: Measuring transferability of adversarial examples (ℓ_∞ , $\epsilon = 8/255$). We use a ten step PGD attack in top figure and FGSM attack in bottom figure. Lower robust accuracy implies better transferability.

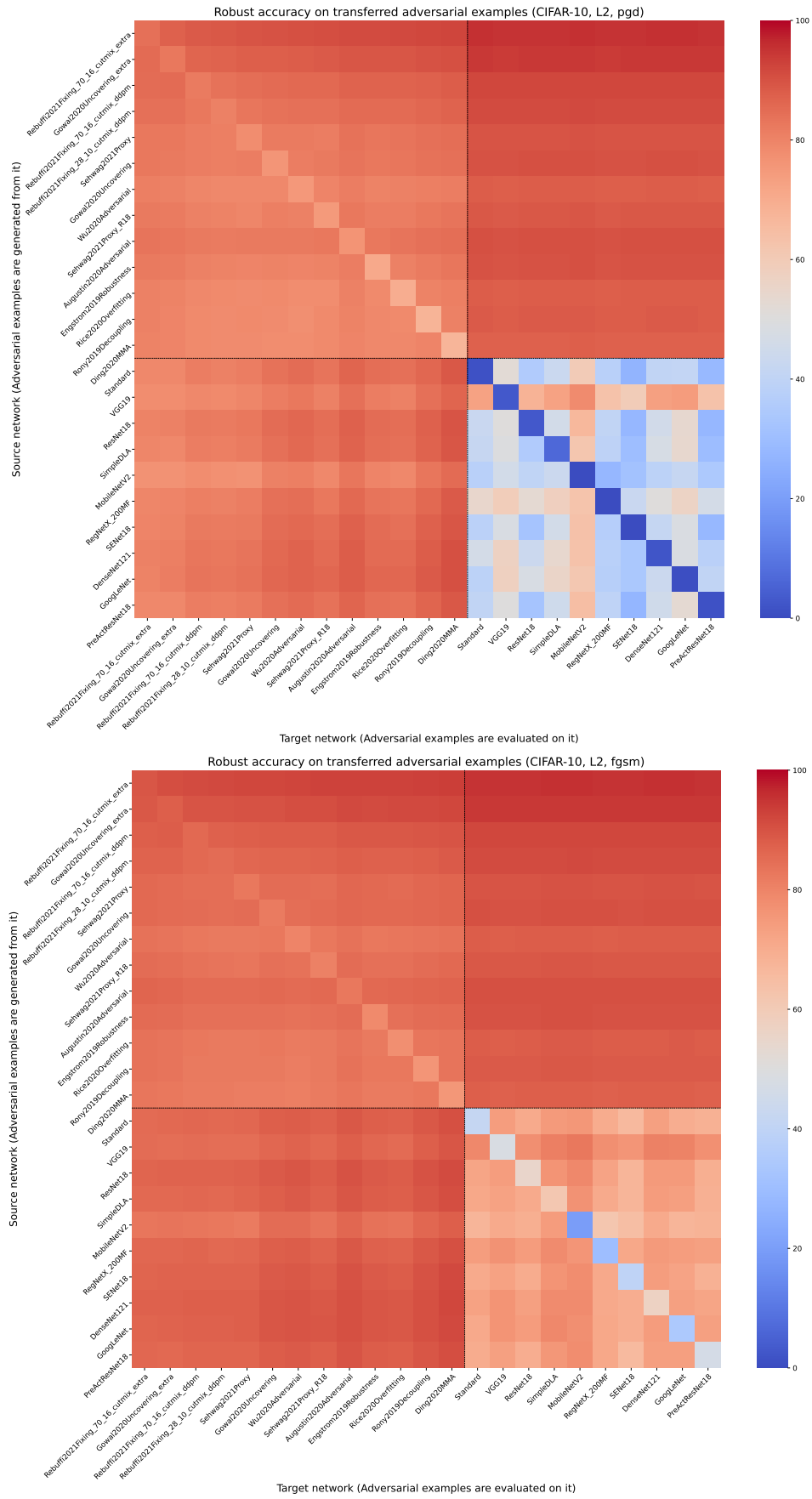


Figure 19: Measuring transferability of adversarial examples (ℓ_2 , $\epsilon = 0.5$). We use a ten step PGD attack in top figure and FGSM attack in bottom figure. Lower robust accuracy implies better transferability.

Table 4: Leaderboard for the ℓ_2 -threat model, CIFAR-10.

	Model	Clean	AA	Extra data	Architecture	Venue	Model Zoo ID
1	Rebuffi et al. [100]	95.74	82.32	Y	WRN-70-16	arXiv, Mar 2021	Rebuffi2021Fixing_70_16_cutmix_extra
2	Gowal et al. [46]	94.74	80.53	Y	WRN-70-16	arXiv, Oct 2020	Gowal2020Uncovering_extra
3	Rebuffi et al. [100]	92.41	80.42	N	WRN-70-16	arXiv, Mar 2021	Rebuffi2021Fixing_70_16_cutmix_ddpm
4	Rebuffi et al. [100]	91.79	78.80	N	WRN-28-10	arXiv, Mar 2021	Rebuffi2021Fixing_28_10_cutmix_ddpm
5	Schwag et al. [112]	90.31	76.12	N	WRN-34-10	arXiv, Apr 2021	Schwag2021Proxy
6	Gowal et al. [46]	90.90	74.50	N	WRN-70-16	arXiv, Oct 2020	Gowal2020Uncovering
7	Wu et al. [140]	88.51	73.66	N	WRN-34-10	NeurIPS 2020	Wu2020Adversarial
8	Schwag et al. [112]	89.52	73.39	N	ResNet-18	arXiv, Apr 2021	Schwag2021Proxy_R18
9	Augustin et al. [7]	91.08	72.91	Y	ResNet-50	ECCV 2020	Augustin2020Adversarial
10	Engstrom et al. [35]	90.83	69.24	N	ResNet-50	GitHub, Sep 2019	Engstrom2019Robustness
11	Rice et al. [102]	88.67	67.68	N	ResNet-18	ICML 2020	Rice2020Overfitting
12	Rony et al. [103]	89.05	66.44	N	WRN-28-10	CVPR 2019	Rony2019Decoupling
13	Ding et al. [32]	88.02	66.09	N	WRN-28-4	ICLR 2020	Ding2020MMA
14	Standard	94.78	0.0	N	WRN-28-10	N/A	Standard

Table 5: Leaderboard for common corruptions, CIFAR-10.

	Model	Clean	Corr.	Extra data	Architecture	Venue	Model Zoo ID
1	Calian et al. [14]	94.93	92.17	Y	ResNet-50	arXiv, Apr 2021	N/A
2	Hendrycks et al. [56]	95.83	89.09	N	ResNeXt29_32x4d	ICLR 2020	Hendrycks2020AugMix_ResNeXt
3	Hendrycks et al. [56]	95.08	88.82	N	WRN-40-2	ICLR 2020	Hendrycks2020AugMix_WRN
4	Kireev et al. [67]	94.77	88.53	N	ResNet-18	arXiv, Mar 2021	Kireev2021Effectiveness_RLATAugMixNoJSD
5	Gowal et al. [46]	94.74	87.68	Y	WRN-70-16	arXiv, Oct 2020	N/A
6	Kireev et al. [67]	94.97	86.60	N	ResNet-18	arXiv, Mar 2021	Kireev2021Effectiveness_AugMixNoJSD
7	Kireev et al. [67]	93.24	85.04	N	ResNet-18	arXiv, Mar 2021	Kireev2021Effectiveness_Gauss50percent
8	Gowal et al. [46]	90.90	84.90	N	WRN-70-16	arXiv, Oct 2020	N/A
9	Kireev et al. [67]	93.10	84.10	N	ResNet-18	arXiv, Mar 2021	Kireev2021Effectiveness_RLAT
10	Gowal et al. [46]	91.10	81.84	Y	WRN-70-16	arXiv, Oct 2020	N/A
11	Gowal et al. [46]	85.29	76.37	N	WRN-70-16	arXiv, Oct 2020	N/A
12	Standard	94.78	73.46	N	WRN-28-10	N/A	Standard

Table 6: Leaderboard for the ℓ_∞ -threat model, CIFAR-100.

	Model	Clean	AA	Extra data	Architecture	Venue	Model Zoo ID
1	Gowal et al. [46]	69.15	36.88	Y	WRN-70-16	arXiv, Oct 2020	Gowal2020Uncovering_extra
2	Rebuffi et al. [100]	63.56	34.64	N	WRN-70-16	arXiv, Mar 2021	Rebuffi2021Fixing_70_16_cutmix_ddpm
3	Rebuffi et al. [100]	62.41	32.06	N	WRN-28-10	arXiv, Mar 2021	Rebuffi2021Fixing_28_10_cutmix_ddpm
4	Cui et al. [28]	62.55	30.20	N	WRN-34-20	arXiv, Nov 2020	Cui2020Learnable_34_20_LBGAT6
5	Gowal et al. [46]	60.86	30.03	N	WRN-70-16	arXiv, Oct 2020	Gowal2020Uncovering
6	Cui et al. [28]	60.64	29.33	N	WRN-34-10	arXiv, Nov 2020	Cui2020Learnable_34_10_LBGAT6
7	Wu et al. [140]	60.38	28.86	N	WRN-34-10	NeurIPS 2020	Wu2020Adversarial
8	Hendrycks et al. [55]	59.23	28.42	Y	WRN-28-10	ICML 2019	Hendrycks2019Using
9	Cui et al. [28]	70.25	27.16	N	WRN-34-10	arXiv, Nov 2020	Cui2020Learnable_34_10_LBGAT0
10	Chen et al. [20]	62.15	26.94	N	WRN-34-10	arXiv, Oct 2020	Chen2020Efficient
11	Sitawarin et al. [117]	62.82	24.57	N	WRN-34-10	ICML 2020	Sitawarin2020Improving
12	Rice et al. [102]	53.83	18.95	N	Pre-activation ResNet 18	ICML 2020	Rice2020Overfitting

Table 7: Leaderboard for common corruptions, CIFAR-100.

	Model	Clean	Corr.	Extra data	Architecture	Venue	Model Zoo ID
1	Hendrycks et al. [56]	78.90	65.54	N	ResNeXt29_32x4d	ICLR 2020	Hendrycks2020AugMix_ResNeXt
2	Hendrycks et al. [56]	76.28	64.63	N	WRN-40-2	ICLR 2020	Hendrycks2020AugMix_WRN
3	Gowal et al. [46]	69.15	56.72	Y	WRN-70-16	arXiv, Oct 2020	N/A
4	Gowal et al. [46]	60.86	50.17	N	WRN-70-16	arXiv, Oct 2020	N/A