
OODRobustBench: a Benchmark and Large-Scale Analysis of Adversarial Robustness under Distribution Shift

Lin Li¹ Yifei Wang² Chawin Sitawarin³ Michael Spratling^{1,4}

Abstract

Existing works have made great progress in improving adversarial robustness, but typically test their method only on data from the same distribution as the training data, i.e. in-distribution (ID) testing. As a result, it is unclear how such robustness generalizes under input distribution shifts, i.e. out-of-distribution (OOD) testing. To address this issue we propose a benchmark named OODRobustBench to comprehensively assess OOD adversarial robustness using 23 dataset-wise shifts (i.e. naturalistic shifts in input distribution) and 6 threat-wise shifts (i.e., unforeseen adversarial threat models). OODRobustBench is used to assess 706 robust models using 60.7K adversarial evaluations. This large-scale analysis shows that: 1) adversarial robustness suffers from a severe OOD generalization issue; 2) ID robustness correlates strongly with OOD robustness in a positive linear way. The latter enables the prediction of OOD robustness from ID robustness. We then predict and verify that existing methods are unlikely to achieve high OOD robustness. Novel methods are therefore required to achieve OOD robustness beyond our prediction. To facilitate the development of these methods, we investigate a wide range of techniques and identify several promising directions. Code and models are available at: <https://github.com/OODRobustBench/OODRobustBench>.

1. Introduction

Adversarial attack poses a serious threat to real-world machine learning models, and various approaches have been developed to defend against such attacks. Previous work (Athalye et al., 2018) has shown that adversarial evaluation

¹Department of Informatics, King’s College London, UK ²MIT CSAIL, USA ³UC Berkeley, USA ⁴University of Luxembourg, Luxembourg. Correspondence to: Lin Li <lin.3.li@kcl.ac.uk>.

is critical to the study of adversarial robustness since an unreliable evaluation can often give a false sense of robustness. However, we believe that even state-of-the-art evaluation benchmarks (like RobustBench Croce et al., 2021) suffer from a severe limitation: they only consider ID generalization where test data comes from the same distribution as the training data. Since distribution shifts are inevitable in the real world, it is crucial to assess how adversarial robustness is affected when the test distribution differs from the training one.

Although OOD generalization has been extensively studied for clean accuracy (Hendrycks & Dietterich, 2019; Taori et al., 2020; Miller et al., 2021; Baek et al., 2022; Zhao et al., 2022; Yang et al., 2022), there is little known about the OOD generalization of adversarial robustness. To fill this void, this paper presents for the first time, a comprehensive benchmark, **OODRobustBench**, for assessing out-of-distribution adversarial robustness. Furthermore, it reports results of a large-scale analysis of existing robust models performed using the new benchmark to answer the following questions:

1. How resilient are current adversarially robust models to distribution shift?
2. Can we predict OOD robustness from ID robustness?
3. How can we achieve OOD robustness?

OODRobustBench is analogous and complementary to RobustBench which is used for assessing in-distribution adversarial robustness. It includes two categories of distribution shift: dataset shift and threat shift (see Figure 1). Dataset shift denotes test data that has different characteristics from the training data due to varying conditions under which the samples are collected: for images, these include but are not limited to corruptions, background, and viewpoint. OODRobustBench contains 23 such dataset shifts and assesses adversarial robustness to such data using the attack seen by the model during training. Threat shift denotes a variation between training and test adversarial threat models. In other words, threat shift assesses a model’s robustness to unseen adversarial attacks applied to ID test data. OODRobustBench employs six different types of threat shifts. Adversarial robustness is evaluated for each type of shift to comprehensively assess OOD robustness.

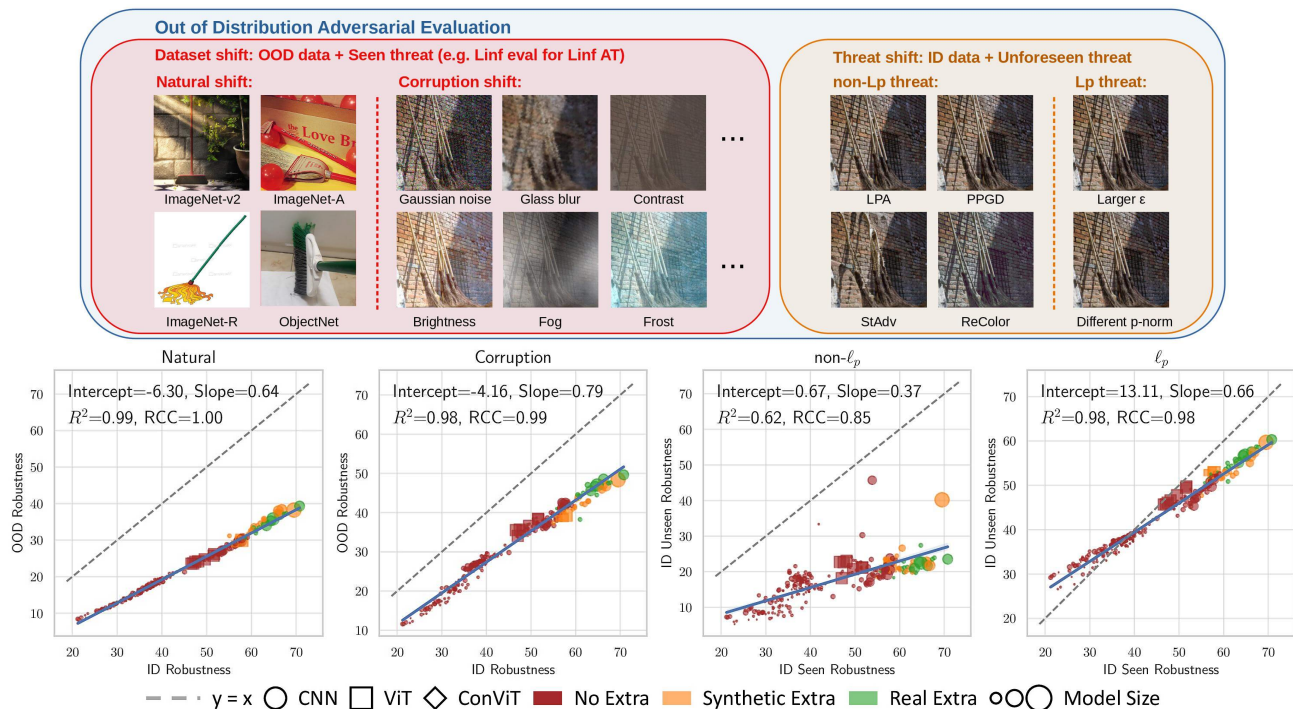


Figure 1. The construction of OODRobustBench (top) and the correlation between ID and OOD robustness under 4 types of distribution shift for CIFAR10 ℓ_∞ (bottom). Each marker represents a model and is annotated by its training set-up. The solid blue line is the fitted linear correlation. The dashed gray line ($y = x$) represents perfect generalization where OOD robustness equals ID robustness. Deviation from the dashed line indicates robustness degradation under the respective distribution shift.

With OODRobustBench, we analyze the OOD generalization behavior of 706 well-trained robust models (a total of 60.7K adversarial evaluations). This model zoo covers a diversity of architectures, robust training methods, data augmentation techniques and training set-ups to ensure the conclusions drawn from this assessment are general and comprehensive. This large-scale analysis reveals that:

- **Adversarial robustness suffers from a severe OOD generalization issue.** Robustness degrades on average by 18%/31%/24% under distribution shifts for CIFAR10 ℓ_∞ , CIFAR10 ℓ_2 and ImageNet ℓ_∞ respectively.
- **ID and OOD accuracy/robustness have a strong linear correlation under many shifts** (visualized in Figure 1). This enables the prediction of OOD performance from ID performance.

The findings above are rigorously identified by a large-scale, systematic, analysis for the first time. Furthermore, our analysis also offer several novel insights into the OOD generalization behavior of adversarial robustness:

- **The higher the ID robustness of the model, the more robustness degrades under distribution shift.** This suggests that while great progress has been made on improving ID robustness, we only gain *diminishing returns* under distribution shift.

- **An abnormal catastrophic drop in robustness under noise shifts is observed in some methods.** For instance, under Gaussian noise shift, HAT (Rade & Moosavi-Dezfooli, 2022) suffers from a severe drop of robustness by 46% whereas the average drop is 9%.
- **Adversarial training boosts the correlation between ID and OOD performance under corruption shifts**, and thus, improves the fidelity of using ID performance for model selection and OOD performance prediction.
- **ℓ_p robustness correlates poorly with non- ℓ_p robustness.** This suggests that non- ℓ_p robustness cannot be predicted from ℓ_p robustness, and enhancing ℓ_p robustness does not necessarily result in improved non- ℓ_p robustness.

Last, we investigate how to achieve OOD adversarial robustness. First, based on the discovered linear trend, we predict the best available OOD performance for the existing ℓ_p -based robustness methodology and find that **existing methods are unlikely to achieve high OOD adversarial robustness** (e.g. the predicted upper bound of OOD robustness under the dataset shifts is only 43% on ImageNet ℓ_∞). Next, we examine a wide range of techniques for achieving OOD adversarial robustness beyond the above prediction. Most of these techniques, including training with extra data, data augmentation, advanced model architectures, scaling-up models and unsupervised representation learning, have

limited or no benefit. However, we do identify several adversarial training methods (Dai et al., 2022; Pang et al., 2020; Ding et al., 2020; Bai et al., 2023) that have the potential to exceed the prediction and produce higher OOD adversarial robustness.

Overall, this work reveals that most existing robust models including the state-of-the-art ones are vulnerable to distribution shifts and demonstrates that the existing approaches to improve ID robustness may be insufficient to achieve high OOD robustness. To ensure safe deployment in the wild, we advocate for the assessment of OOD robustness in future models and for the development of new approaches that can cope with distribution shifts better and achieve OOD robustness beyond our prediction.

2. Related Works

Robustness under dataset shift. Early work (Sehwag et al., 2019) studied the generalization of robustness to novel classes that are unseen during training. On the other hand, our setup only considers the input distribution shift and not the unforeseen classes. Recently, Sun et al. (2022b) studied the OOD generalization of certified robustness under corruption shifts for a few state-of-the-art methods. In contrast, we focus on empirical robustness instead of certified robustness. Alhamoud et al. (2023) is the most relevant work. They studied the generalization of robustness from multiple source domains to an unseen domain. Different from them, the models we examine are trained on only one source domain, which is the most common set-up in the existing works of adversarial training (Croce et al., 2021). Moreover, we also cover much more diverse distribution shifts, models and training methods than Sun et al. (2022b) and Alhamoud et al. (2023) so that the conclusion drawn in this work is more general and comprehensive.

Robustness against unforeseen adversarial threat models. It was observed that naive adversarial training (Madry et al., 2018) with only one single ℓ_p threat model generalizes poorly to unforeseen ℓ_p threat models, e.g., higher perturbation bound (Stutz et al., 2020), different p -norm (Tramer & Boneh, 2019; Maini et al., 2020; Croce & Hein, 2022), or non- ℓ_p threat models including color transformation ReColor (Laidlaw & Feizi, 2019), spatial transformation StAdv (Xiao et al., 2018), LPIPS-bounded attacks PPGD and LPA (Laidlaw et al., 2021) and many others (Kaufmann et al., 2023). We complement the existing works by conducting a large-scale analysis on the unforeseen robustness of ℓ_p robust models trained by varied methods and training set-ups. We are thus able to provide new insights into the generalization of robustness to unforeseen threat models and identify effective yet previously unknown approaches to enhance unforeseen robustness.

More related works are discussed in Appendix A.

3. OOD Adversarial Robustness Benchmark

3.1. OODRobustBench

OODRobustBench is designed to simulate the possible data distribution shifts that might occur in the wild and evaluate adversarial robustness in the face of them. It focuses on two types of distribution shifts: dataset shift and threat shift. *Dataset shift*, OOD_d , denotes the distributional difference between training and test raw datasets. *Threat shift*, OOD_t , denotes the difference between training and evaluation *threat models*, a special type of distribution shift. The original test set drawn from the same distribution as the training set is considered ID. The variant dataset with the same classes yet where the distribution of the inputs differs is considered OOD.

Dataset shift. To represent diverse data distribution in the wild, OODRobustBench includes multiple types of dataset shifts from two sources: *natural* and *corruption*. For natural shifts, we adopt four different variant datasets per source dataset: CIFAR10.1 (Recht et al., 2018), CIFAR10.2 (Lu et al., 2020), CINIC (Darlow et al., 2018), and CIFAR10-R (Hendrycks et al., 2021a) for CIFAR10, and ImageNet-v2 (Recht et al., 2019), ImageNet-A (Hendrycks et al., 2021b), ImageNet-R (Hendrycks et al., 2021a), and ObjectNet (Barbu et al., 2019) for ImageNet. For corruption shifts, we adopt, from the corruption benchmarks (Hendrycks & Dietterich, 2019), 15 types of common corruption in four categories: Noise (gaussian, impulse, shot), Blur (motion, defocus, glass, zoom), Weather (fog, snow, frost) and Digital (brightness, contrast, elastic, pixelate, JPEG). Each corruption has five levels of severity. Overall, the dataset-shift testbed consists of 79 ($4 + 15 \times 5$) subsets. Appendix B.1 gives the details of the above datasets and data processing.

Accuracy and robustness are evaluated on the ID and OOD dataset. To compute the overall performance of OOD_d , we first average the result of natural and corruption shifts:

$$R_c(f) = \mathbb{E}_{i \in \{\text{corruptions}\}, j \in \{\text{severity}\}} R_{i,j}(f) \quad (1)$$

$$R_n(f) = \mathbb{E}_{i \in \{\text{naturals}\}} R_i(f) \quad (2)$$

where $R(\cdot)$ returns accuracy or adversarial robustness and f denotes the model to be assessed. Next, we average the above two results to get the overall performance of the dataset shift as

$$R_{ood}(f) = (R_c(f) + R_n(f))/2 \quad (3)$$

To evaluate a model, OODRobustBench performs 80 (79 for OOD_d and 1 for ID) runs of adversarial evaluation. This makes computationally expensive attacks like AutoAttack (Croce & Hein, 2020) impractical to use. To balance efficiency and effectiveness, we use MM5 (Gao et al., 2022) for

Table 1. Performance, evaluated with OODRobustBench, of state-of-the-art models trained on CIFAR10 to be robust to ℓ_∞ attacks. The top 3 results for each metric are highlighted in **bold** and/or underscore. Significant ranking discrepancies are indicated in **red**. The "OOD" column presents the average of the robustness to OOD_d and OOD_t . The complete leaderboard, featuring a total of 396 models, is available at <https://oodrobustbench.github.io/>.

Method	Model	Accuracy (%)		Robustness (%)				Ranking	
		ID	OOD_d	ID	OOD_d	OOD_t	OOD	ID	OOD
BDM (Wang et al., 2023)	WRN-70-16	93.2	76.0	70.7	44.4	35.8	40.1	1	2
AS (Bai et al., 2023)	WRN-70-16 + ResNet-152	95.2	79.0	69.5	43.3	46.7	45.0	2	1
DKL (Cui et al., 2023)	WRN-28-10	92.1	74.8	<u>67.7</u>	42.4	35.4	38.9	3	4
BDM (Wang et al., 2023)	WRN-28-10	92.4	75.0	<u>67.3</u>	42.3	35.2	38.8	4	5
FDA (Rebuffi et al., 2021)	WRN-70-16	92.2	74.8	66.7	<u>42.6</u>	33.6	38.1	5	9
DDPM (Gowal et al., 2021b)	WRN-70-16	88.7	70.6	66.2	42.7	33.6	38.2	6	8
Uncovering (Gowal et al., 2021a)	WRN-70-16	91.1	73.2	66.0	42.5	34.0	38.2	7	7
RobustResNet (Huang et al., 2023a)	WRN-A4	91.5	73.8	65.8	41.7	33.3	37.5	8	12
FDA (Rebuffi et al., 2021)	WRN-106-16	88.5	70.6	64.8	41.4	33.9	37.6	9	10
DyART (Xu et al., 2022)	WRN-28-10	93.6	77.2	64.7	39.6	<u>37.0</u>	38.3	10	6
PORT (Schwag et al., 2022)	ResNet-152	87.2	69.2	62.7	40.7	<u>32.3</u>	36.5	17	15
HAT (Rade & Moosavi-Dezfooli, 2022)	WRN-28-10	88.1	69.4	60.9	35.1	30.2	32.6	22	57
AWP (Wu et al., 2020)	WRN-28-10	88.2	69.8	60.1	38.2	31.3	34.8	26	27
RST (Carmon et al., 2019)	WRN-28-10	89.6	71.5	59.8	36.7	31.1	33.9	28	38
MART (Wang et al., 2020)	WRN-28-10	87.5	70.2	56.7	35.5	32.6	34.0	52	35
HE (Pang et al., 2020)	WRN-34-20	85.1	66.9	53.8	32.4	46.2	39.3	70	3
FAT (Zhang et al., 2020)	WRN-34-10	84.5	65.9	53.6	32.9	31.8	32.4	71	59
Overfitting (Rice et al., 2020)	WRN-34-20	85.3	66.4	53.5	32.0	27.8	29.9	72	89
TRADES (Zhang et al., 2019)	WRN-34-10	84.9	66.5	52.6	31.6	26.5	29.1	76	99
FBF (Wong et al., 2020)	PreActResNet-18	83.3	64.9	43.3	25.3	24.8	25.0	111	112

robustness evaluation. MM5 is approximately $32\times$ faster than AutoAttack (Gao et al., 2022) while achieving similar results, as verified in Appendix B.2 alongside the results of evaluations using alternative attacks. The perturbation bound ϵ is $8/255$ for CIFAR10 ℓ_∞ , 0.5 for CIFAR10 ℓ_2 and $4/255$ for ImageNet ℓ_∞ .

Threat shift. OODRobustBench adopts six unforeseen attacks as in Laidlaw et al. (2021); Dai et al. (2022) to simulate threat shifts. They are categorized into two groups, ℓ_p and non- ℓ_p , according to whether they are bounded by the ℓ_p norm or not. The ℓ_p shift group includes MM attacks with the same p -norm but larger ϵ and with different p -norm. The non- ℓ_p shift group includes the imperceptible, PPGD and LPA, and perceptible, ReColor and StAdv, attacks. The overall robustness under threat shift, OOD_t , is simply the mean of these six unforeseen attacks. These attacks are selected because they cover a wide range of different scenarios of threat shift and each of them is representative of its corresponding category (100+ cites). We are aware of alternative non- ℓ_p attacks (Kaufmann et al., 2023) but do not include them due to the constraint of computational resource.

We follow the same setting as Laidlaw et al. (2021); Dai et al. (2022) to configure the above attacks since this has been well tested to be effective. The ℓ_p attacks use $\epsilon = 12/255$ and $\epsilon = 0.5$ for ℓ_∞ and ℓ_2 threats on CIFAR10 ℓ_∞ , $\epsilon = 8/255$ and $\epsilon = 1$ for ℓ_∞ and ℓ_2 threats on CIFAR10 ℓ_2 and on ImageNet ℓ_∞ . The perturbation bound is 0.5 for PPGD, 0.5 for LPA, 0.05 for StAdv and 0.06 for ReColor. The number of iterations is 40 for PPGD and LPA regardless of dataset, is 100 for StAdv and ReColor on CIFAR10 and 200 on ImageNet.

Criteria for robust models are described in Appendix C.1 and are the same as RobustBench (Croce et al., 2021).

3.2. OOD Performance and Ranking

The benchmark results for CIFAR10 ℓ_∞ , ℓ_2 and ImageNet ℓ_∞ are in Tables 1, 4 and 5 respectively.

Robustness degrades significantly under distribution shift. For models trained to be robust for CIFAR10 ℓ_∞ (Figure 2), CIFAR10 ℓ_2 (Figure 9) and ImageNet ℓ_∞ (Figure 10), the average drop in robustness (ID adversarial accuracy - OOD adversarial accuracy) is $18\%/20\%/27\%$ under dataset shift and $18\%/42\%/22\%$ under threat shift. Robustness degradation is much severe for a subset of shifts: whereas the average robustness degradation of OOD_d is 18% on CIFAR10 ℓ_∞ , some shifts like CIFAR10-R, fog and contrast degrade by 38% , 30% and 32% , respectively.

The higher the ID robustness of the model, the more robustness degrades under the shifts. For example, the top method in Table 1 degrades by 30% of robustness, while the bottom method degrades by only 18% . This suggests that while the great progress has been made on improving ID robustness, we only gain diminishing returns under the distribution shifts. Besides, in Figure 2, the distribution of robustness degradation for most shifts spreads over a wide range, suggesting a large variation across individual models.

Robustness degradation under noise shifts can be abnormally catastrophic (the outliers under noise shifts in Figure 2). This issue is most severe on (Rade & Moosavi-Dezfooli, 2022) whose robustness falls by $43\%/46\%/38\%$ under impulse/Gaussian/shot noise, whereas the average drop is $12\%/9\%/8\%$ (discussed in Appendix E). A similar

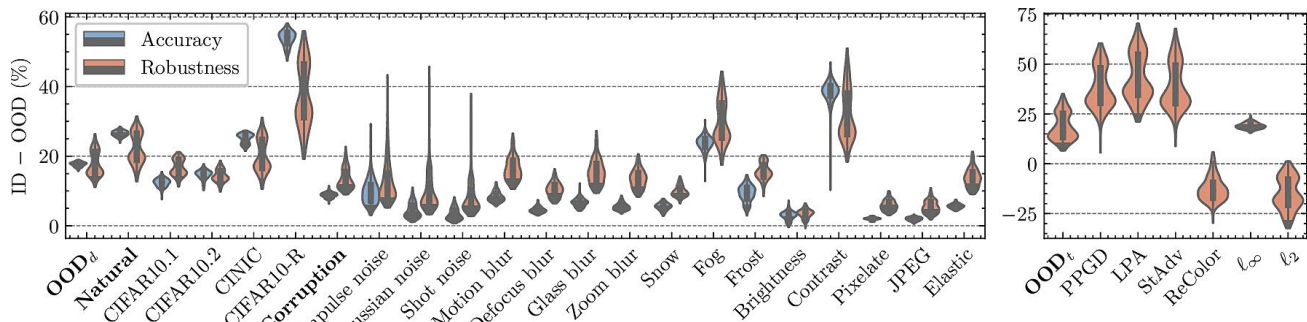


Figure 2. Degradation of accuracy and robustness under various distribution shifts for CIFAR10 ℓ_∞ .

yet milder drop is also observed on [Debenedetti et al. \(2023\)](#) and models trained with some advanced data augmentations like AutoAugment ([Cubuk et al., 2019](#)).

Higher ID robustness generally implies higher OOD robustness but not always (see the last two columns of Tables 1, 4 and 5). For example, in Table 1, the ranking of [Rade & Moosavi-Dezfooli \(2022\)](#) drops from 22 to 57 due to catastrophic degradation, while the ranking of [Pang et al. \(2020\)](#) jumps from 70 to 3 due to its superior robustness under threat shift (analyzed in Section 6.3).

4. Linear Trend and OOD Prediction

It was previously observed that OOD accuracy is strongly correlated with ID accuracy under many dataset shifts for Standardly-Trained (ST) models ([Miller et al., 2021](#)). This property is important since it enables the model selection and OOD performance prediction through ID performance. Nevertheless, it is unclear if such correlation still holds for adversarial robustness. This is particularly intriguing because accuracy and robustness usually go in opposite directions: i.e. there is a trade-off between accuracy and robustness ([Tsipras et al., 2019](#)). Furthermore, the threat shifts as a scenario of OOD are unique to adversarial evaluation and were, thus, never explored in the previous studies of accuracy trends. Surprisingly, we find that ID and OOD robustness also have a linear correlation under many distribution shifts. It is even more surprising that the correlation for AT models is much stronger than that for ST models.

The following result is based on a large-scale analysis including over 60K OOD evaluations of 706 models. 187 of these models were retrieved from RobustBench or other published works so as to include current state-of-the-art methods, and the remaining models were trained by ourselves. These models are mainly trained in three set-ups: CIFAR10 ℓ_∞ , CIFAR10 ℓ_2 and ImageNet ℓ_∞ . They cover a wide range of model architectures, model sizes, data augmentation methods, training and regularization techniques. More detail is given in Appendix C.2.

4.1. Linear Trend under Dataset Shift

This section studies how ID and OOD accuracy/robustness correlate under dataset shifts. We fit a linear regression on four pairs of metrics (Acc-Acc, Rob-Rob, Acc-Rob, and Rob-Acc) for each dataset shift and each training setup (CIFAR10 ℓ_∞ , CIFAR10 ℓ_2 and ImageNet ℓ_∞). Taking Acc-Rob as an example, a linear model is fitted with ID accuracy as the observed variable x and OOD adversarial robustness as the target variable y . The result for each shift is given in Appendix H. Below are the major findings.

ID accuracy (resp. robustness) strongly correlates with OOD accuracy (resp. robustness) in a linear relationship for most dataset shifts. In Figures 3, 11 and 12, the regression of Acc-Acc and Rob-Rob for most shifts achieve very high R^2 (> 0.9), i.e., their relationship can be well explained by a linear model. This suggests for these shifts ID performance is a good indication of OOD performance, and more importantly, OOD performance can be reliably predicted by ID performance using the fitted linear model.

Nevertheless, under some shifts, ID and OOD performance are only weakly correlated. Natural shifts like CIFAR10-R and ImageNet-A and corruption shifts like noise, fog and contrast are observed to have relatively low R^2 across varied training set-ups in Figures 3, 11 and 12. It can be seen from Figures 16 and 17 that the correlation for these shifts becomes even weaker, and the gap of R^2 between them and the others expands, as more inferior (relatively worse accuracy and/or robustness) models are excluded from the regression. This suggests that the models violating the linear trend are mostly high-performance. Appendix F discusses how inferior models are identified and how they influence the correlation.

AT models exhibit a stronger linear correlation between ID and OOD accuracy under most corruption shifts on CIFAR10 in Figures 3 and 11. The improvement is dramatic for particular shifts. For example, R^2 surges from nearly 0 (no linear correlation) for ST models to around 0.8 (evident linear correlation) for AT models with Gaussian and shot noise data shifts. These results are contrary to the previ-

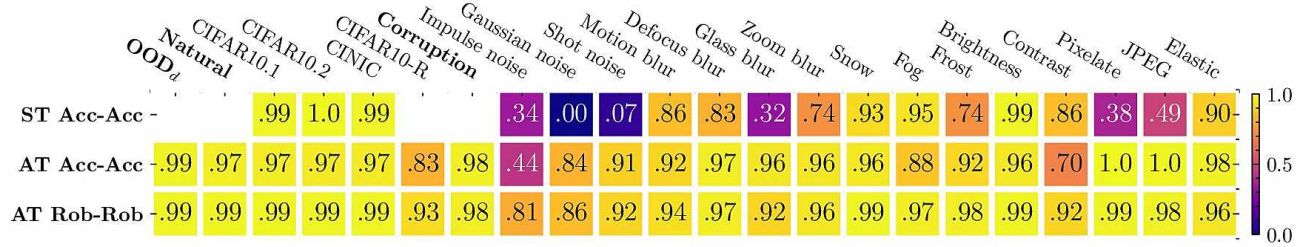


Figure 3. R^2 of regression between ID and OOD performance for Standardly-Trained (ST) and Adversarially-Trained (AT) models under various dataset shifts for CIFAR10 ℓ_∞ . Higher R^2 implies stronger linear correlation. The results for ST models were copied from (Miller et al., 2021). Some results of ST are missing (blank cells) because they were not reported in (Miller et al., 2021).

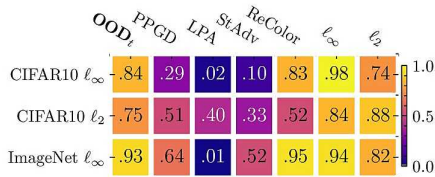


Figure 4. R^2 of regression between ID seen robustness and OOD unforeseen robustness, i.e., threat shift.

ous finding on ST models (Miller et al., 2021). However, note that we measure linear correlation for the raw data, whereas (Miller et al., 2021) applies a nonlinear transformation to the data to promote linearity. Overall, adversarial training boosts linear correlation for corruption shifts, and hence, improves the faithfulness of using ID performance for model selection and OOD performance prediction.

We attribute this to AT improving accuracy on the corrupted data (Kireev et al., 2022). Intuitively, ST models have less correlated corruption accuracy because corruption significantly impairs accuracy and such effect varies a lot among models. Compared to ST, AT effectively mitigates the effect of corruption on accuracy, and hence, reduces the divergence of corruption accuracy so that corruption accuracy is more correlated to ID accuracy.

Last, we observe **no evident correlation when ID and OOD metrics misalign, i.e., Acc-Rob and Rob-Acc for CIFAR10**, but weak correlation for ImageNet ℓ_∞ as shown in Figure 13. This is due to the varied trade-off between accuracy and robustness of different models (discussed in details in Appendix F.1)

4.2. Linear Trend under Threat Shift

This section studies the relationship between seen and unforeseen robustness. Both seen and unforeseen robustness are computed using only ID data yet with different attacks. Linear regression is then conducted between seen robustness (x) and unforeseen robustness (y). The result of regression for each threat shift is given in Appendix I. The sensitivity of the regression results to the composition of the model zoo is discussed in Appendix F.

ℓ_p robustness correlates poorly with non- ℓ_p robustness. R^2 of the regression between ID ℓ_p robustness and PPGD,

LPA and StAdv robustness is low in Figure 4. Particularly, R^2 is close to 0 for ℓ_∞ -LPA and ℓ_∞ -StAdv on CIFAR10 ℓ_∞ suggesting no correlation at all. As shown in Figures 28 to 30, the increase in ID ℓ_p robustness leads to only slight or even no improvement on unforeseen robustness esp. for LPA and StAdv. Interestingly, despite poor correlation with PPGD, LPA and StAdv, ID ℓ_p robustness is well correlated with ReColor unforeseen robustness.

ℓ_p robustness correlates strongly with ℓ_p robustness of different ϵ and p -norm. R^2 of their regression is higher than 0.7 across all assessed set-ups in Figure 4 suggesting a consistently strong linear correlation. The correlation between different ϵ of the same p -norm is stronger than the correlation between different p -norm.

4.3. Unsupervised OOD Robustness Prediction

The linear trends discovered above enable the prediction of OOD performance only if labeled OOD data is available. There is a line of works (Baek et al., 2022; Deng & Zheng, 2021; Garg et al., 2021) showing that OOD accuracy can be predicted with only unlabeled OOD data. We study here if OOD adversarial robustness can be predicted, similarly, in an unsupervised manner. We run the experiments with CIFAR-10 ℓ_∞ models for CIFAR-10.1 (Figure 5) and Impulse noise (Figure 14) shifts and find that a linear trend is also observed in the agreement between the predictions of any pair of two robust models: R^2 is 0.99 for CIFAR-10.1 shift and 0.95 for Impulse noise shift. This suggests that the unsupervised method (Baek et al., 2022) is also effective in predicting OOD adversarial robustness.

5. Incompetence in OOD Generalization

Based on the precise linear trend observed above for existing robust training methods, we can predict the OOD performance of a model trained by such a method from its ID performance using the fitted linear model. Furthermore, we can extrapolate from current trends to predict the maximum OOD robustness that can be expected from a hypothetical future model that achieves perfect robustness on ID data (assuming the linear trend continues).

$$\text{slope} \times 100 + \text{intercept}. \quad (4)$$

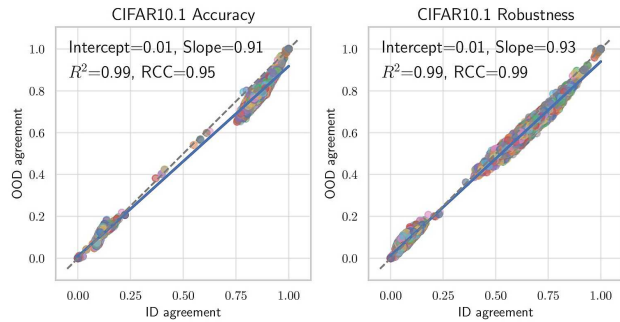


Figure 5. Correlation between ID and OOD prediction agreement on adversarial examples for CIFAR10 ℓ_∞ AT models. Each point represents the prediction agreement of two models.

This estimates the best OOD performance one can expect by fully exploiting existing robust training techniques. Note that a wide range of models and techniques (Appendix C.2) are covered by our correlation analysis so their, as well as their variants’, OOD performance should be (approximately) bounded by the predicted upper limit. The accuracy of the prediction depends on the R^2 of the correlation.

We find that **continuously improving ID ℓ_p robustness following existing practice is unlikely to achieve high OOD adversarial robustness**. The upper limit of OOD robustness under dataset shift, OOD_d , is 66%/71%/43% for CIFAR10 ℓ_∞ (Figure 6), CIFAR10 ℓ_2 (Figure 15) and ImageNet ℓ_∞ (Figure 15) respectively, and under threat shift OOD_t is 52%/35%/52% correspondingly. Hence, if current trends continue, the resulting models are likely to be very unreliable in real-world applications. The vulnerability of models is most evident for ImageNet ℓ_∞ under dataset shift and for CIFAR10 ℓ_2 under threat shift. The expected upper limit of OOD robustness also varies greatly among individual shifts ranging from nearly 0 to 100%.

One of the accounts for this issue is that the existing methods have poor conversion rate to OOD robustness from ID robustness as shown by the slope of the linear trend in Figures 6 and 15. Taking an example of fog shift on ImageNet, the slope is roughly 0.1 so improving 10% ID robustness can only lead to 1% improvement on fog robustness. Besides, the upper limit and conversion rate of robustness are observed to be much lower than those of accuracy in Figure 15, suggesting the OOD generalization issue is more severe for robustness. Overall, this issue calls for developing novel methods that can improve OOD robustness beyond our prediction.

6. Improving OOD Adversarial Robustness

To inspire the design of methods that have OOD robustness exceeding the above prediction, this section investigates methods that have the potential to be effective for boosting the OOD generalization of robustness. The effectiveness is

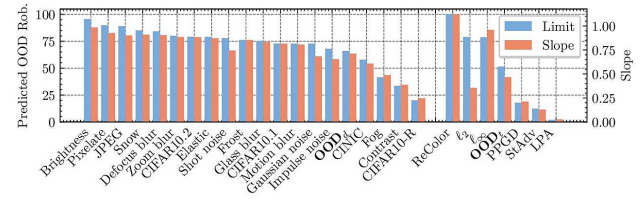


Figure 6. The estimated upper limit of OOD robustness, and the slope, of OOD robustness from ID robustness under various distribution shifts for CIFAR10 ℓ_∞ . The estimated upper limit is capped by 100% as robustness can not surpass 100%.

quantified by two metrics: OOD performance and effective performance. Effective performance measures the extra resilience of a model under distribution shift when compared to a group of models by adapting the metric of “Effective Robustness” (Taori et al., 2020):

$$R'(f) = R_{ood}(f) - \beta(R_{id}(f)) \quad (5)$$

where $\beta(\cdot)$ is a linear mapping from ID to OOD metric fitted on a group of models. We name this metric effective robustness (adversarial effective robustness) when R_{id} and R_{ood} are accuracy (robustness). A positive adversarial effective robustness means that f achieves adversarial robustness above what the linear trend predicts based on its ID performance, i.e., f is advantageous over the fitted models on OOD generalization. Note that higher adversarial effective robustness is not equivalent to higher OOD robustness since the model may have a lower ID robustness. The specific set-ups and detailed results of the following experiments are described in Appendix G.

6.1. Data

Training with extra data boosts both robustness and adversarial effective robustness compared to training schemes without extra data (see Figure 7a). There is no clear advantage to training with extra real data (Carmon et al., 2019) rather than synthetic data (Gowal et al., 2021b) except for the adversarial effective robustness under threat shift which is improved more by real data.

Advanced data augmentation improves robustness under both types of shifts and adversarial effective robustness under threat shift over the baseline augmentation RandomCrop (see Figure 7b). Nevertheless, advanced data augmentation methods other than TA (Müller & Hutter, 2021) degrade adversarial effective robustness under dataset shift.

6.2. Model

Advanced model architecture greatly boosts robustness and adversarial effective robustness under both types of shift over the baseline ResNet (He et al., 2016) (Figure 7c). Among all tested architectures, ViT (Dosovitskiy et al., 2021) achieves the highest adversarial effective robustness.

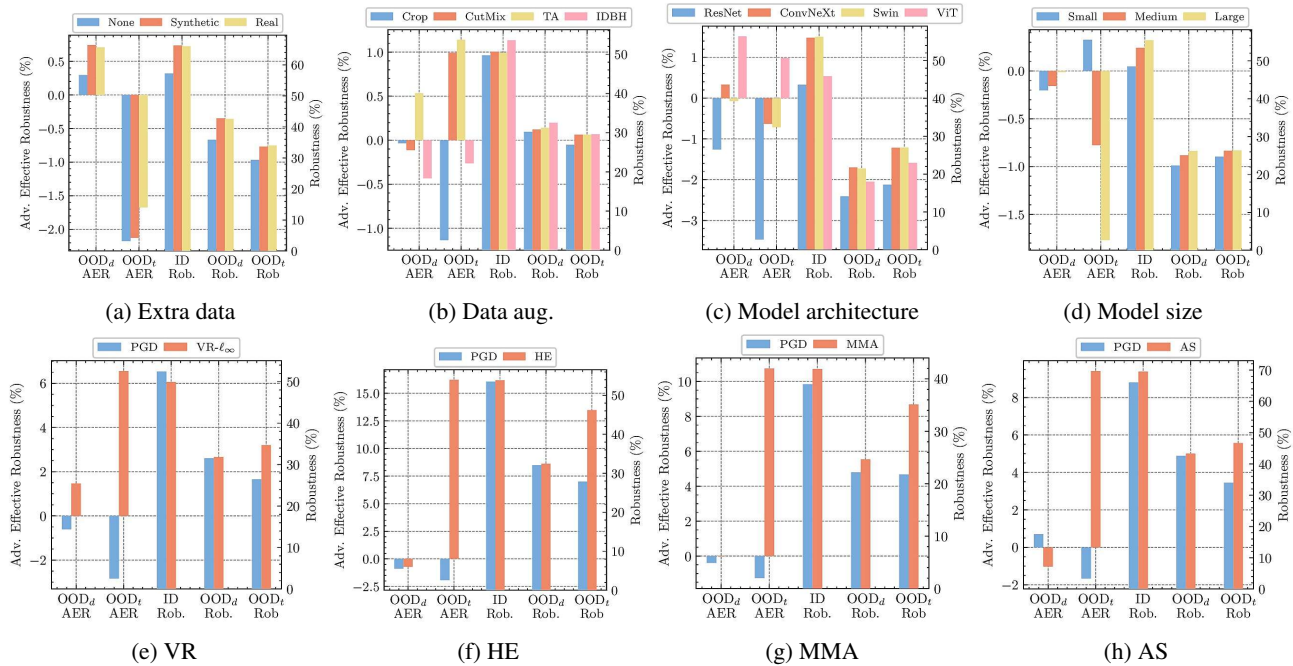


Figure 7. The robustness (Rob.) and adversarial effective robustness (AER) of various robust techniques.

Scaling model up improves robustness under both types of shift and adversarial effective robustness under dataset shift, but dramatically impairs adversarial effective robustness under threat shift (Figure 7d). The latter is because increasing model size greatly improves ID robustness but not OOD robustness so that the real OOD robustness is much below the OOD robustness predicted by linear correlation.

6.3. Adversarial Training

VR (Dai et al., 2022), the state-of-the-art defense against unforeseen attacks, greatly boosts adversarial effective robustness under threat shifts in spite of inferior ID robustness. Surprisingly, VR also clearly boosts adversarial effective robustness under dataset shift even though not designed for dealing with these shifts.

Training methods **HS** (Pang et al., 2020), **MMA** (Ding et al., 2020) and **AS** (Bai et al., 2023) achieve an AER of 16.22%, 10.74% and 9.41%, respectively, under threat shift, which are much higher than the models trained with PGD. Importantly, in contrast to VR, these methods also improve ID robustness resulting in a further boost on OOD robustness. This makes them a potentially promising defense against multi-attack (Dai et al., 2023).

6.4. OOD Generalization Methods

Two leading methods, CARD-Deck (Diffenderfer et al., 2021) (ranked 1st) and PLAT (Kireev et al., 2022), from the common corruptions leaderboard of RobustBench are evaluated using our benchmark in Table 2. Despite the expected

remarkable OOD clean generalization under OOD_d shifts, they offer little or no adversarial robustness regardless of ID or OOD setting. It suggests that OOD generalization methods alone do not help OOD adversarial robustness unless combined with adversarial training.

To test the effect of combining adversarial training with OOD generalization method, we evaluate a recent attempt in this direction, MSD+REx (Ibrahim et al., 2023). This approach trains using the multi-attack defense MSD with various attacks and applies REx (Krueger et al., 2021) by treating different attacks as separate domains. Surprisingly, as shown in Table 2, this purpose-built solution impairs OOD adversarial robustness under both dataset and threat shifts and offers no evident improvement in AER when compared to supervised ℓ_p adversarial training.

While these findings suggest that this specific implementation might be ineffective, the combination of AT and OOD methods remains a promising direction. Future work should focus on a more careful design for this integration. One potential strategy could be to treat different groups of data, rather than attacks, as different domains like Huang et al. (2023b).

6.5. Unsupervised Representation Learning

Unsupervised learning has been observed to train models that generalize to distribution shifts better than supervised learning (Shi et al., 2023; Shen et al., 2021). However, it is unclear whether or not unsupervised learning will benefit OOD adversarial robustness. To test this we evaluated a

Table 2. The performance of non-single- ℓ_p defense and OOD generalization methods under distribution shift on CIFAR10 ℓ_∞ . The AER of StAdv AT, PLAT and CARD-Deck is invalid (“-”) because of their (nearly) 0 ID/OOD robustness.

Defense	Method	Model	ID		OOD _d				OOD _t		
			Acc.	Rob.	Acc.	Rob.	ER	AER	non- ℓ_p	ℓ_p	Avg.
Supervised ℓ_p defense	PGD ℓ_p AT	ResNet18	83.4	49.4	64.6	<u>30.2</u>	-1.1	0.2	19.1	45.5	27.9
Self-supervised ℓ_p defense	ACL	ResNet18	82.3	49.4	64.2	30.1	-0.6	0.1	20.2	<u>45.1</u>	28.5
non- ℓ_p attack defense	ReColor AT	ResNet50	93.4	8.5	78.0	3.4	2.8	2.6	16.4	18.1	17.0
	StAdv AT	ResNet50	86.2	0.1	66.8	0.0	-1.6	-	9.5	3.4	7.4
	PAT-Alexnet	ResNet50	71.6	28.6	57.1	16.8	2.3	1.6	<u>48.4</u>	33.1	<u>43.3</u>
	PAT-Self	ResNet50	82.4	30.4	66.3	16.8	1.4	0.3	32.1	36.6	33.6
Unforeseen attack defense	VR	ResNet18	72.9	<u>48.9</u>	56.4	31.4	0.5	<u>1.7</u>	24.8	43.3	31.0
	PAT+VR	ResNet50	72.5	<u>29.4</u>	56.8	17.4	1.3	<u>1.7</u>	55.0	33.6	47.8
Composite attack defense	GAT-f	ResNet50	82.3	38.7	66.2	22.2	1.4	-0.2	14.8	17.0	15.5
	GAT-fs	ResNet50	82.1	41.9	65.7	24.8	1.2	0.1	17.2	18.3	17.5
Multi-attack defense	MAAT-Average	ResNet50	86.8	39.9	70.7	22.4	1.7	-0.8	25.0	41.7	30.5
	MAAT-Max	ResNet50	84.0	25.6	68.1	13.1	1.8	0.0	30.2	29.8	30.0
	MAAT-Random	ResNet50	85.2	22.1	67.7	10.5	0.2	0.0	12.7	31.0	18.8
OOD generalization method	PLAT	ResNet18	94.7	0.1	80.3	0.0	4.0	-	0.0	0.0	0.0
	CARD-Deck	WRN-18-2	96.5	1.0	83.5	0.5	5.4	-	0.0	0.0	0.0
Multi-attack + OOD defense	MSD+REx	ResNet18	78.0	43.3	59.9	26.2	-0.7	0.5	19.0	40.4	26.1

model trained by Adversarial Contrastive Learning (ACL) (Jiang et al., 2020) which combines self-supervised contrastive learning with adversarial training. The effective robustness under dataset shift is 0.1% (Table 2), suggesting only marginal benefit in improving OOD robustness.

6.6. Non-single- ℓ_p Defenses

This section evaluates the OOD generalization capability of various defenses beyond the supervised single-attack ℓ_p defense. We compared several specific methods, including AT with the color-based attack ReColor (Laidlaw et al., 2021), the spatial attack StAdv (Laidlaw et al., 2021), and the LPIPS-bound attack PAT-Alexnet/Self (Laidlaw et al., 2021). Additionally, we examined composite attacks defense (e.g., color plus ℓ_p) GAT-f/fs (Hsiung et al., 2023) and multiple attacks adversarial training (MAAT) (Tramer & Boneh, 2019; Maini et al., 2020) involving ℓ_2 , ℓ_∞ , StAdv, and ReColor. MAAT has three variants: "Average" optimizes the average loss across all attacks, "Max" optimizes the maximum loss across all attacks, and "Random" selects a random attack at each training iteration.

Unfortunately, none of these defenses achieve high OOD_d ER and AER in Table 2, indicating that they are not significantly better than the supervised single-attack ℓ_p AT at handling OOD dataset distribution shifts. This reinforces our conclusion that achieving OOD adversarial robustness is challenging with existing methods and underscores the need to develop new approaches that effectively address distribution shifts.

While MAAT and PAT show substantial improvements over ℓ_p AT in terms of robustness against non- ℓ_p attacks, this is partly because some of the non- ℓ_p attacks used were already

encountered during their training, making them no longer unforeseen. This highlights the difficulty in benchmarking unforeseen robustness across different types of defenses.

6.7. Summary

The evaluated techniques, except for some AT methods (Section 6.3), achieve relatively limited or even no adversarial effective robustness. This suggests that applying them is unlikely to significantly change the linear trend in Section 4 and thus the predicted upper limit of OOD robustness (Section 5). In contrast, the methods identified in Section 6.3 show the promise in achieving OOD performance beyond our prediction. Another promising direction is to combine OOD generalization methods with adversarial training.

7. Conclusions

This work proposes a new benchmark to assess OOD adversarial robustness, provides many insights into the generalization of existing robust models under distribution shift and identifies several robust interventions beneficial to OOD generalization. We have analyzed the OOD robustness of hundreds of diverse models to ensure that we obtain generally applicable insights. As we focus on general trends, our analysis does not provide a detailed investigation into individual methods or explain the observed outliers such as the catastrophic robustness degradation. However, OODRobustBench provides a tool for performing such more detailed investigations in the future. It also provides a means of measuring progress towards models that are more robust in real-world conditions and will, hopefully, spur the future development of such models.

Acknowledgment

The authors acknowledge the use of the research computing facility at King’s College London, King’s Computational Research, Engineering and Technology Environment (CRE-ATE). Lin Li was funded by the King’s - China Scholarship Council (K-CSC). Yifei Wang was supported by Office of Naval Research under grant N00014-20-1-2023 (MURI ML-SCOPE), NSF AI Institute TILOS (NSF CCF-2112665), and NSF award 2134108. Chawin was supported in part by funds provided by the National Science Foundation (under grant 2229876), the KACST-UCB Center for Secure Computing, the Department of Homeland Security, IBM, the Noyce Foundation, Google, Open Philanthropy, and the Center for AI Safety Compute Cluster. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the sponsors.

Impact Statement

This paper presents work whose goal is to advance the field of adversarial machine learning. There are many potential societal consequences of our work, none of which we feel are negative and must be specifically highlighted here.

References

- Alhamoud, K., Hammoud, H. A. A. K., Alfarra, M., and Ghanem, B. Generalizability of Adversarial Robustness Under Distribution Shifts. *Transactions on Machine Learning Research*, May 2023.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *Proceedings of the 35th International Conference on Machine Learning*, July 2018.
- Augustin, M., Meinke, A., and Hein, M. Adversarial Robustness on In- and Out-Distribution Improves Explainability. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- Baek, C., Jiang, Y., Raghunathan, A., and Kolter, J. Z. Agreement-on-the-line: Predicting the performance of neural networks under distribution shift. In *Advances in Neural Information Processing Systems*, volume 35, pp. 19274–19289, 2022.
- Bai, Y., Anderson, B. G., Kim, A., and Sojoudi, S. Improving the Accuracy-Robustness Trade-Off of Classifiers via Adaptive Smoothing, May 2023.
- Barbu, A., Mayo, D., Alverio, J., Luo, W., Wang, C., Gutfreund, D., Tenenbaum, J., and Katz, B. ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In *Advances in Neural Information Processing Systems*, 2019.
- Carmon, Y., Raghunathan, A., Schmidt, L., Duchi, J. C., and Liang, P. S. Unlabeled Data Improves Adversarial Robustness. In *Advances in Neural Information Processing Systems*, 2019.
- Croce, F. and Hein, M. Reliable Evaluation of Adversarial Robustness with an Ensemble of Diverse Parameter-free Attacks. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- Croce, F. and Hein, M. Adversarial robustness against multiple and single l_p -threat models via quick fine-tuning of robust classifiers. In *International Conference on Machine Learning*, pp. 4436–4454. PMLR, 2022.
- Croce, F., Andriushchenko, M., Sehwag, V., Debenedetti, E., Flammarion, N., Chiang, M., Mittal, P., and Hein, M. RobustBench: a standardized adversarial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, October 2021.
- Croce, F., Goyal, S., Brunner, T., Shelhamer, E., Hein, M., and Cemgil, T. Evaluating the adversarial robustness of adaptive test-time defenses. In *International Conference on Machine Learning*, pp. 4421–4435. PMLR, 2022.
- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. AutoAugment: Learning Augmentation Strategies From Data. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Cui, J., Tian, Z., Zhong, Z., Qi, X., Yu, B., and Zhang, H. Decoupled Kullback-Leibler Divergence Loss, May 2023.
- Dai, S., Mahloujifar, S., and Mittal, P. Formulating Robustness Against Unforeseen Attacks. In *Advances in Neural Information Processing Systems*, December 2022.
- Dai, S., Mahloujifar, S., Xiang, C., Sehwag, V., Chen, P.-Y., and Mittal, P. MultiRobustBench: Benchmarking Robustness Against Multiple Attacks. In *International Conference on Machine Learning*, May 2023.
- Darlow, L. N., Crowley, E. J., Antoniou, A., and Storkey, A. J. CINIC-10 is not ImageNet or CIFAR-10, October 2018.
- Debenedetti, E., Sehwag, V., and Mittal, P. A Light Recipe to Train Robust Vision Transformers. In *First IEEE Conference on Secure and Trustworthy Machine Learning*, February 2023.

- Deng, W. and Zheng, L. Are Labels Always Necessary for Classifier Accuracy Evaluation? In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021.
- Diffenderfer, J., Bartoldson, B., Chaganti, S., Zhang, J., and Kailkhura, B. A Winning Hand: Compressing Deep Networks Can Improve Out-of-Distribution Robustness. In *Advances in Neural Information Processing Systems*, 2021.
- Ding, G. W., Sharma, Y., Lui, K. Y. C., and Huang, R. MMA Training: Direct Input Space Margin Maximization through Adversarial Training. In *International Conference on Learning Representations*, 2020.
- Dong, Y., Fu, Q.-A., Yang, X., Pang, T., Su, H., Xiao, Z., and Zhu, J. Benchmarking Adversarial Robustness on Image Classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- Dong, Y., Ruan, S., Su, H., Kang, C., Wei, X., and Zhu, J. ViewFool: Evaluating the Robustness of Visual Recognition to Adversarial Viewpoints. In *Advances in Neural Information Processing Systems*, December 2022.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Hounsby, N. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*, 2021.
- Engstrom, L., Ilyas, A., Salman, H., Santurkar, S., and Tsipras, D. Robustness (python library), 2019. URL <https://github.com/MadryLab/robustness>.
- Ford, N., Gilmer, J., Carlini, N., and Cubuk, D. Adversarial examples are a natural consequence of test error in noise. 2019. doi: 10.48550/arXiv.1901.10513.
- Gao, R., Wang, J., Zhou, K., Liu, F., Xie, B., Niu, G., Han, B., and Cheng, J. Fast and Reliable Evaluation of Adversarial Robustness with Minimum-Margin Attack. In *Proceedings of the 39th International Conference on Machine Learning*, June 2022.
- Garg, S., Balakrishnan, S., Lipton, Z. C., Neyshabur, B., and Sedghi, H. Leveraging unlabeled data to predict out-of-distribution performance. October 2021.
- Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019.
- Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann, F. A. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–73, 2020. doi: 10.1038/s42256-020-00257-z.
- Gowal, S., Qin, C., Uesato, J., Mann, T., and Kohli, P. Uncovering the Limits of Adversarial Training against Norm-Bounded Adversarial Examples. *arXiv:2010.03593 [cs, stat]*, March 2021a.
- Gowal, S., Rebuffi, S.-A., Wiles, O., Stimberg, F., Calian, D., and Mann, T. Improving Robustness using Generated Data. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021b.
- Guo, C., Rana, M., Cisse, M., and Maaten, L. v. d. Countering Adversarial Images using Input Transformations. In *International Conference on Learning Representations*, 2018.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Hendrycks, D. and Dietterich, T. Benchmarking Neural Network Robustness to Common Corruptions and Perturbations. In *International Conference on Learning Representations*, 2019.
- Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F., Dorundo, E., Desai, R., Zhu, T., Parajuli, S., Guo, M., Song, D., Steinhardt, J., and Gilmer, J. The Many Faces of Robustness: A Critical Analysis of Out-of-Distribution Generalization. In *International Conference on Computer Vision*, 2021a.
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural Adversarial Examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021b.
- Hsiung, L., Tsai, Y.-Y., Chen, P.-Y., and Ho, T.-Y. Towards Compositional Adversarial Robustness: Generalizing Adversarial Training to Composite Semantic Perturbations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- Hu, J., Shen, L., and Sun, G. Squeeze-and-Excitation Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely Connected Convolutional Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

- Huang, S., Lu, Z., Deb, K., and Boddeti, V. N. Revisiting Residual Networks for Adversarial Robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023a.
- Huang, Z., Zhu, M., Xia, X., Shen, L., Yu, J., Gong, C., Han, B., Du, B., and Liu, T. Robust generalization against photon-limited corruptions via worst-case sharpness minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16175–16185, 2023b.
- Ibrahim, A., Guille-Escuret, C., Mitliagkas, I., Rish, I., Krueger, D., and Bashivan, P. Towards Out-of-Distribution Adversarial Robustness, February 2023.
- Jiang, Z., Chen, T., Chen, T., and Wang, Z. Robust Pre-Training by Adversarial Contrastive Learning. In *Advances in Neural Information Processing Systems*, 2020.
- Kaufmann, M., Kang, D., Sun, Y., Basart, S., Yin, X., Mazeika, M., Arora, A., Dziedzic, A., Boenisch, F., Brown, T., Steinhardt, J., and Hendrycks, D. Testing Robustness Against Unforeseen Adversaries, July 2023.
- Kireev, K., Andriushchenko, M., and Flammarion, N. On the effectiveness of adversarial training against common corruptions. In *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, August 2022.
- Krueger, D., Caballero, E., Jacobsen, J.-H., Zhang, A., Binas, J., Zhang, D., Le Priol, R., and Courville, A. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, pp. 5815–5826. PMLR, 2021.
- Laidlaw, C. and Feizi, S. Functional Adversarial Attacks. In *Advances in Neural Information Processing Systems*, 2019.
- Laidlaw, C., Singla, S., and Feizi, S. Perceptual Adversarial Robustness: Defense Against Unseen Threat Models. In *International Conference on Learning Representations*, January 2021.
- Li, L. and Spratling, M. Improved Adversarial Training Through Adaptive Instance-wise Loss Smoothing, March 2023a.
- Li, L. and Spratling, M. Understanding and combating robust overfitting via input loss landscape analysis and regularization. *Pattern Recognition*, April 2023b.
- Li, L. and Spratling, M. W. Data augmentation alone can improve adversarial training. In *International Conference on Learning Representations*, February 2023c.
- Li, L., Qiu, J., and Spratling, M. AROID: Improving Adversarial Robustness through Online Instance-wise Data Augmentation, June 2023.
- Liu, C., Dong, Y., Xiang, W., Yang, X., Su, H., Zhu, J., Chen, Y., He, Y., Xue, H., and Zheng, S. A Comprehensive Study on Robustness of Image Classification Models: Benchmarking and Rethinking, February 2023.
- Lu, S., Nott, B., Olson, A., Todeschini, A., Vahabi, H., Carmon, Y., and Schmidt, L. Harder or Different? A Closer Look at Distribution Shift in Dataset Reproduction. In *ICML 2020 Workshop on Uncertainty and Robustness in Deep Learning*, 2020.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*, 2018.
- Maini, P., Wong, E., and Kolter, Z. Adversarial Robustness Against the Union of Multiple Perturbation Models. In *Proceedings of the 37th International Conference on Machine Learning*, November 2020.
- Mao, X., Chen, Y., Li, X., Qi, G., Duan, R., Zhang, R., and Xue, H. Easyrobust: A comprehensive and easy-to-use toolkit for robust computer vision, 2022.
- Miller, J. P., Taori, R., Raghunathan, A., Sagawa, S., Koh, P. W., Shankar, V., Liang, P., Carmon, Y., and Schmidt, L. Accuracy on the Line: on the Strong Correlation Between Out-of-Distribution and In-Distribution Generalization. In *Proceedings of the 38th International Conference on Machine Learning*, July 2021.
- Müller, S. G. and Hutter, F. TrivialAugment: Tuning-Free Yet State-of-the-Art Data Augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.
- Pang, T., Yang, X., Dong, Y., Xu, K., Zhu, J., and Su, H. Boosting Adversarial Training with Hypersphere Embedding. In *Advances in Neural Information Processing Systems*, 2020.
- Pang, T., Lin, M., Yang, X., Zhu, J., and Yan, S. Robustness and Accuracy Could Be Reconcilable by (Proper) Definition. In *Proceedings of the 39th International Conference on Machine Learning*, June 2022.
- Rade, R. and Moosavi-Dezfooli, S.-M. Reducing Excessive Margin to Achieve a Better Accuracy vs. Robustness Trade-off. In *International Conference on Learning Representations*, March 2022.
- Rebuffi, S.-A., Goyal, S., Calian, D. A., Stimberg, F., Wiles, O., and Mann, T. Fixing Data Augmentation to Improve Adversarial Robustness. Technical report, October 2021.

- Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do CIFAR-10 Classifiers Generalize to CIFAR-10? *arXiv:1806.00451 [cs, stat]*, June 2018.
- Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do ImageNet Classifiers Generalize to ImageNet? In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Rice, L., Wong, E., and Kolter, J. Z. Overfitting in adversarially robust deep learning. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- Rony, J., Hafemann, L. G., Oliveira, L. S., Ben Ayed, I., Sabourin, R., and Granger, E. Decoupling Direction and Norm for Efficient Gradient-Based L2 Adversarial Attacks and Defenses. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Rusak, E., Schott, L., Zimmermann, R. S., Bitterwolf, J., Bringmann, O., Bethge, M., and Brendel, W. A simple way to make neural networks robust against diverse image corruptions. In *Proceedings of the European Conference on Computer Vision*, 2020.
- Samangouei, P., Kabkab, M., and Chellappa, R. DefenseGAN: Protecting Classifiers Against Adversarial Attacks Using Generative Models. In *International Conference on Learning Representations*, February 2018.
- Sandler, M., Howard, A. G., Zhu, M., Zhmoginov, A., and Chen, L.-C. Mobilenetv2: Inverted residuals and linear bottlenecks. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018. doi: 10.1109/CVPR.2018.00474.
- Sehwag, V., Bhagoji, A. N., Song, L., Sitawarin, C., Cullina, D., Chiang, M., and Mittal, P. Analyzing the Robustness of Open-World Machine Learning. In *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security*, November 2019.
- Sehwag, V., Mahloujifar, S., Handina, T., Dai, S., Xiang, C., Chiang, M., and Mittal, P. Robust Learning Meets Generative Models: Can Proxy Distributions Improve Adversarial Robustness? In *International Conference on Learning Representations*, March 2022.
- Shen, Z., Liu, J., He, Y., Zhang, X., Xu, R., Yu, H., and Cui, P. Towards Out-Of-Distribution Generalization: A Survey, August 2021.
- Shi, Y., Daunhawer, I., Vogt, J. E., Torr, P. H., and Sanyal, A. How robust is unsupervised representation learning to distribution shift? In *The Eleventh International Conference on Learning Representations*. OpenReview, 2023.
- Simonyan, K. and Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *International Conference on Learning Representations*, 2015.
- Singh, N. D., Croce, F., and Hein, M. Revisiting Adversarial Training for ImageNet: Architectures, Training and Generalization across Threat Models. 2023.
- Stutz, D., Hein, M., and Schiele, B. Confidence-Calibrated Adversarial Training: Generalizing to Unseen Attacks. In *International Conference on Machine Learning*, 2020.
- Sun, J., Mehra, A., Kailkhura, B., Chen, P.-Y., Hendrycks, D., Hamm, J., and Mao, Z. M. Certified adversarial defenses meet out-of-distribution corruptions: Benchmarking robustness and simple baselines. In *Proceedings of the European Conference on Computer Vision*, 2022a.
- Sun, J., Mehra, A., Kailkhura, B., Chen, P.-Y., Hendrycks, D., Hamm, J., and Mao, Z. M. A Spectral View of Randomized Smoothing Under Common Corruptions: Benchmarking and Improving Certified Robustness. In *Computer Vision – ECCV 2022*, 2022b.
- Szegedy, C., Wei Liu, Yangqing Jia, Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. Going deeper with convolutions. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the Inception Architecture for Computer Vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Tan, M. and Le, Q. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. In *Proceedings of the 36th International Conference on Machine Learning*, May 2019.
- Tang, S., Gong, R., Wang, Y., Liu, A., Wang, J., Chen, X., Yu, F., Liu, X., Song, D., Yuille, A., Torr, P. H. S., and Tao, D. RobustART: Benchmarking Robustness on Architecture Design and Training Techniques, January 2022.
- Taori, R., Dave, A., Shankar, V., Carlini, N., Recht, B., and Schmidt, L. Measuring Robustness to Natural Distribution Shifts in Image Classification. In *Advances in Neural Information Processing Systems*, 2020.
- Torralba, A., Fergus, R., and Freeman, W. T. 80 Million Tiny Images: A Large Data Set for Nonparametric Object and Scene Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, November 2008.
- Tramer, F. and Boneh, D. Adversarial Training and Robustness for Multiple Perturbations. In *Advances in Neural Information Processing Systems*, 2019.

- Trockman, A. and Kolter, J. Z. Patches are all you need? *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. Featured Certification.
- Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., and Madry, A. Robustness May Be at Odds with Accuracy. In *International Conference on Learning Representations*, 2019.
- Wang, H., Ge, S., Lipton, Z., and Xing, E. P. Learning Robust Global Representations by Penalizing Local Predictive Power. In *Advances in Neural Information Processing Systems*, 2019.
- Wang, Y., Zou, D., Yi, J., Bailey, J., Ma, X., and Gu, Q. Improving Adversarial Robustness Requires Revisiting Misclassified Examples. In *International Conference on Learning Representations*, 2020.
- Wang, Z., Pang, T., Du, C., Lin, M., Liu, W., and Yan, S. Better Diffusion Models Further Improve Adversarial Training. In *International Conference on Machine Learning*, February 2023.
- Wong, E., Rice, L., and Kolter, J. Z. Fast is better than free: Revisiting adversarial training. In *International Conference on Learning Representations*, 2020.
- Wu, D., Xia, S.-T., and Wang, Y. Adversarial Weight Perturbation Helps Robust Generalization. In *Advances in Neural Information Processing Systems*, 2020.
- Xiao, C., Zhu, J.-Y., Li, B., He, W., Liu, M., and Song, D. Spatially Transformed Adversarial Examples. February 2018.
- Xie, S., Girshick, R., Dollar, P., Tu, Z., and He, K. Aggregated Residual Transformations for Deep Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- Xu, Y., Sun, Y., Goldblum, M., Goldstein, T., and Huang, F. Exploring and Exploiting Decision Boundary Dynamics for Adversarial Robustness. September 2022.
- Yang, J., Wang, P., Zou, D., Zhou, Z., Ding, K., Peng, W., Wang, H., Chen, G., Li, B., Sun, Y., Du, X., Zhou, K., Zhang, W., Hendrycks, D., Li, Y., and Liu, Z. OpenOOD: Benchmarking generalized out-of-distribution detection, October 2022.
- Yu, F., Wang, D., Shelhamer, E., and Darrell, T. Deep layer aggregation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2403–2412, 2018. doi: 10.1109/CVPR.2018.00255.
- Zhang, H., Yu, Y., Jiao, J., Xing, E., Ghaoui, L. E., and Jordan, M. Theoretically Principled Trade-off between Robustness and Accuracy. In *International Conference on Machine Learning*, May 2019.
- Zhang, J., Xu, X., Han, B., Niu, G., Cui, L., Sugiyama, M., and Kankanhalli, M. Attacks Which Do Not Kill Training Make Adversarial Learning Stronger. In *Proceedings of the 37th International Conference on Machine Learning*, November 2020.
- Zhang, R., Isola, P., Efros, A. A., Shechtman, E., and Wang, O. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.
- Zhao, B., Yu, S., Ma, W., Yu, M., Mei, S., Wang, A., He, J., Yuille, A., and Kortylewski, A. OOD-CV: A benchmark for robustness to out-of-distribution shifts of individual nuisances in natural images. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VIII*, pp. 163–180, Berlin, Heidelberg, 2022. Springer-Verlag. ISBN 978-3-031-20073-1. doi: 10.1007/978-3-031-20074-8_10.

A. Additional Related Works

Except for a few exceptions (Geirhos et al., 2020; Sun et al., 2022a; Rusak et al., 2020; Ford et al., 2019), previous work on generalization to input distribution shifts has not considered adversarial robustness. Hence, work on robustness to OOD data and adversarial attacks has generally happened in parallel, as exemplified by RobustBench (Croce et al., 2021) which provides independent benchmarks for assessing performance on corrupt data and adversarial threats.

A line of works (Tramer & Boneh, 2019; Maini et al., 2020) defends against a union of ℓ_p threat models by training with multiple ℓ_p threat models jointly, which makes these threat models no longer unforeseen. PAT (Laidlaw et al., 2021) replaces ℓ_p bound with LPIPS (Zhang et al., 2018) in adversarial training and achieves high robustness against several unforeseen attacks. Alternatively, (Dai et al., 2022) proposes variation regularization in addition to ℓ_p adversarial training and improves unforeseen robustness.

Robustness benchmarks. There is a line of works on benchmarking adversarial robustness, including Dong et al. (2020), RobustBench (Croce et al., 2021), RobustART (Tang et al., 2022), MultiRobustBench (Dai et al., 2023), UA (Kaufmann et al., 2023), and Liu et al. (2023). RobustBench focuses on ID adversarial evaluation. Dong et al. (2020) evaluates adversarial robustness under ℓ_p threat shifts in addition to ID adversarial evaluation. RobustART, compared to Dong et al. (2020), also supports the evaluation of OOD clean accuracy under dataset shifts. MultiRobustBench and UA both extend the evaluation set to include more unforeseen attacks, assessing adversarial robustness under threat shifts. Liu et al. (2023) is similar to RobustART but supports a larger number of OOD datasets. Compared to these benchmarks, OODRobustBench is unique in supporting the evaluation of OOD adversarial robustness under dataset shifts while also incorporating the functionalities of the other benchmarks.

A.1. Comparison with Related Works

Is the linear trend of robustness really expected given the linear trend of accuracy?

No. There is a well-known trade-off between accuracy and robustness in the ID setting (Tsipras et al., 2019). We further confirm this fact for the models we evaluate in Figure 13 in the appendix. This means that accuracy and robustness usually go in opposite directions making the linear trend we discover in both particularly interesting. Furthermore, the threat shifts as a scenario of OOD are unique to adversarial evaluation and were thus never explored in the previous studies of accuracy trends.

How does the linear trends observed by us differ from the previously discovered ones?

Robust models exhibit a stronger linear correlation between ID and OOD accuracy for most corruption shifts (Figure 3). Particularly, the boost on linearity is dramatic for shifts including Impulse, Shot and Gaussian noises, Glass blur, Pixelate, and JPEG. For instance, R2 surges from 0 (no linear correlation) for non-robust models to 0.84 (evident linear correlation) for robust models with Gaussian noise data shifts. This suggests that, for robust models, predicting OOD accuracy from ID accuracy is more faithful and applicable to more shifts.

The linear trend of robustness is even stronger than that of accuracy for dataset shifts (Figure 3) but with a lower slope (Section 5). The latter leads to a predicted upper limit of OOD robustness that is way lower than that of OOD accuracy suggesting that the OOD generalization of robustness is much more challenging.

How does our analysis differ from the similar analysis in the prior works?

The scale of these previous works is rather small. For instance, RobustBench observes linear correlation only for three shifts on CIFAR-10 based on 39 models with either ResNet or WideResNet architectures. In such a narrow setting, it is actually neither surprising to see a linear trend nor reliable for predicting OOD performance. By contrast, our conclusion is derived from much more shifts on CIFAR-10 and ImageNet based on 706 models. Importantly, our model zoo covers a diverse set of architectures, robust training methods, data augmentation techniques, and training set-ups. This makes our conclusion more generalizable and the observed (almost perfect) linear trend much more significant.

Similarly, the existing works only test a few models under threat shifts. Those methods are usually just the baseline AT method plus different architectures or the relevant defenses, e.g., jointly trained with multiple threats. It is unclear how

the state-of-the-art robust models perform under threat shifts. By conducting a large-scale analysis, we find that those SOTA models generalize poorly to other threats while also discovering several methods that have relatively inferior ID performance but superior OOD robustness under threat shift. Our analysis therefore facilitates future works in this direction by identifying what techniques are ineffective and what are promising.

How does your benchmark differ from RobustBench?

Our benchmark focuses on OOD adversarial robustness while RobustBench focuses on ID adversarial robustness. Specifically, our benchmark contrasts RobustBench in the datasets and the attacks. We use CIFAR-10.1, CIFAR-10.2, CINIC, and CIFAR-10-R (ImageNet-V2, ImageNet-A, ImageNet-R, ObjectNet) to simulate input data distribution shift for the source datasets CIFAR-10 (ImageNet), while RobustBench only uses the latter source datasets. We use PPGD, LPA, ReColor, StAdv, Linf-12/255, L2-0.5 (PPGD, LPA, ReColor, StAdv, Linf-8/255, L2-1) to simulate threat shift for the training threats Linf-8/255 (L2-0.5), while RobustBench only evaluates the same threats as the training ones.

B. Benchmark Set-up

B.1. Datasets

This section introduces the OOD datasets of natural shifts. For ImageNet, we have:

- **ImageNet-V2** is a reproduction of ImageNet using a completely new set of images. It has the same 1000 classes as ImageNet and each class has 10 images so 10K images in total.
- **ImageNet-A** is an adversarially-selected reproduction of ImageNet. The images in this dataset were selected to be those most misclassified by an ensemble of ResNet-50s. It has 200 ImageNet classes and 7.5K images.
- **ImageNet-R** contains various artistic renditions of objects from ImageNet, so there is a domain shift. It has 30K images and 200 ImageNet classes.
- **ObjectNet** is a large real-world dataset for object recognition. It is constructed with controls to randomize background, object rotation and viewpoint. It has 313 classes but only 104 classes compatible with ImageNet classes so we only use this subset. The selected subset includes 17.2K images.

For CIFAR10, we have:

- **CIFAR10.1** is a reproduction of CIFAR10 using a completely new set of images. It has 2K images sampled from the same source as CIFAR10, i.e., 80M TinyImages (Torralba et al., 2008). It has the same number of classes as CIFAR10.
- **CIFAR10.2** is another reproduction of CIFAR10. It has 12K (10k for training and 2k for test) images sampled from the same source as CIFAR10, i.e., 80M TinyImages. It has the same number of classes as CIFAR10. We only use the test set of CIFAR10.2.
- **CINIC** is a downsampled subset of ImageNet with the same image resolution and classes as CIFAR10. Its test set has 90K images in total, of which 20K images are from CIFAR10 and 70K images are from ImageNet. We use only the ImageNet part.
- **CIFAR10-R** is a new dataset created by us. The images in CIFAR10-R and CIFAR10 have different styles so there is a domain shift. We follow the same procedure as CINIC to downscale the images from ImageNet-R to the same resolution as CIFAR10 and select images from the classes of ImageNet corresponding to CIFAR10 classes. We follow the same class mapping between ImageNet and CIFAR10 as CINIC. Note that ImageNet-R does not have images of the ImageNet classes corresponding to CIFAR10 classes of "airplane" and "horse", so there are only 8 classes in CIFAR10-R.

In practice, we evaluate models using a random sample of 5K images from each of the ImageNet variant datasets, and 10K images from each of the CIFAR10 variant datasets, if those datasets contain more images than that number. This is done to accelerate the evaluation and follows the practice used in RobustBench (Croce et al., 2021).

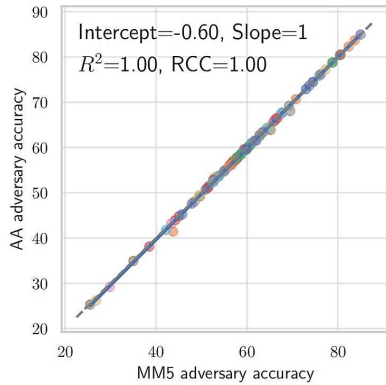


Figure 8. Comparison of MM5 adversarial accuracy against AutoAttack adversarial accuracy. Each data point represents a model.

Table 3. The accuracy and adversarial robustness evaluated by various attacks of models on CIFAR10 and ImageNet for ℓ_∞ threat model. These models are sourced from RobustBench respective leaderboards and identified by their RobustBench identifiers. The results for corruption shifts are reported for severity level 3 .

Dataset	Model	ID					OOD _d				
		Acc.	PGD	CW	AA	MM5	Acc.	PGD	CW	AA	MM5
CIFAR10	Wong2020Fast	83.3	46.6	46.3	43.2	43.3	65.0	27.6	27.3	24.8	25.0
	Engstrom2019Robustness	87.0	52.3	52.6	49.3	49.7	70.6	31.7	31.8	29.1	29.3
	Huang2020Self	83.5	56.2	54.0	52.9	53.0	65.0	34.4	32.9	31.7	31.8
	Sehwag2021Proxy_R18	84.6	58.7	57.2	55.6	55.7	67.0	37.8	36.5	34.6	34.8
	Wang2020Improving	87.5	62.6	58.7	56.3	56.8	70.5	40.3	37.1	34.7	35.1
ImageNet	Salman2020Do_R18	52.9	29.8	27.3	25.3	25.5	21.4	9.5	8.5	7.5	7.7
	Wong2020Fast	55.6	30.0	28.9	26.3	26.8	21.6	8.7	8.5	7.2	7.5
	Engstrom2019Robustness	62.5	32.9	32.6	29.2	29.8	27.2	10.1	10.2	8.4	8.7
	Salman2020Do_R50	64.1	39.0	37.6	34.7	35.0	27.6	12.5	12.0	10.5	10.7
	Singh2023Revisiting_ViT-S-ConvStem	72.6	51.4	50.6	48.1	48.5	39.8	19.8	19.4	17.5	17.7

The OODRobustBench framework is designed for easy integration of alternative datasets representing input data distribution shifts. We plan to maintain and update our code to continually incorporate new OOD datasets such as ImageNet-V (Dong et al., 2022), Stylized-ImageNet (Geirhos et al., 2019), and ImageNet-Sketch (Wang et al., 2019). For the latest developments, we invite readers to visit our GitHub repository: <https://github.com/OODRobustBench/OODRobustBench>.

B.2. Verification of the Effectiveness of the MM5 Attack

B.2.1. COMPARISON OF MM5 AGAINST AUTOATTACK

To verify the effectiveness of MM5, we compare its result with the result of AutoAttack on the ID dataset across all publicly available models from RobustBench for CIFAR10 ℓ_∞ , CIFAR10 ℓ_2 and ImageNet ℓ_∞ . As shown in Figure 8, almost all models¹ are approximately on the line of $y = x$ (gray dashed line) suggesting that their MM5 adversarial accuracy is very close to AA adversarial accuracy. Specifically, the mean gap between MM5 and AA adversarial accuracy is 0.16 and the standard deviation is 0.32.

B.2.2. COMPARISON OF MM5 AGAINST DIVERSE ATTACKS

To test the effectiveness of MM5 against various attacks, we selected 10 models from RobustBench and evaluated their robustness against PGD100, CW100, and AutoAttack. As shown in Table 3, MM5 achieved an adversarial accuracy comparable to that of the strongest attack, AutoAttack, in both in-distribution (ID) and out-of-distribution (OOD) settings. This suggests that the MM5 attack is a reliable method for adversarial evaluation, even under distribution shifts.

¹Two models, Ding et al. (2020) and Xu et al. (2022), are observed to have a slightly higher adversarial accuracy compared to the corresponding AutoAttack results. We use MM+ (Gao et al., 2022) attack to evaluate these two models for a more reliable evaluation and the result of MM+ is close to AutoAttack.

C. Model Zoo

C.1. Criteria for Robust Models

We follow the same criteria as the popular benchmarks (RobustBench (Croce et al., 2021), MultiRobustBench (Dai et al., 2023), etc), which only include robust models that (1) have in general non-zero gradients w.r.t. the inputs, (2) have a fully deterministic forward pass (i.e. no randomness) and (3) do not have an optimization loop. These criteria include most AT models, while excluding most preprocessing methods because they rely on randomness like Guo et al. (2018) or inner optimization loop like Samangouei et al. (2018) which leads to false security, i.e., high robustness to the non-adaptive attack but vulnerable to the adaptive attack.

Meanwhile, we acknowledge that evaluating dynamic preprocessing-based defenses is still an active area of research. It is tricky (Croce et al., 2022), and there has not been a consensus on how to evaluate them. So now, we exclude them for a more reliable evaluation. We will keep maintaining this benchmark, and we would be happy to include them in the future if the community has reached a consensus on that (e.g., if these models are merged into RobustBench).

C.2. Model Zoo

Our model zoo consists of 706 models, of which:

- 396 models are trained on CIFAR10 by ℓ_∞ 8/255
- 239 models are trained on CIFAR10 by ℓ_2 0.5
- 56 models are trained on ImageNet by ℓ_∞ 4/255
- 10 models are trained on CIFAR10 for non- ℓ_p adversarial robustness
- 5 models are trained on CIFAR10 for common corruption robustness

Among the above models, 66 models of CIFAR10 ℓ_∞ , 19 models of CIFAR10 ℓ_2 and 18 models of ImageNet ℓ_∞ are retrieved from RobustBench. 84 models are retrieved from the published works including (Li et al., 2023; Li & Spratling, 2023c;b;a; Liu et al., 2023; Singh et al., 2023; Dai et al., 2022; Hsiung et al., 2023; Mao et al., 2022). The remaining models are trained by ourselves.

We locally train additional models with varying architectures and training parameters to complement the public models from RobustBench on CIFAR-10. We consider 20 model architectures: DenseNet-121 (Huang et al., 2017), GoogLeNet (Szegedy et al., 2015), Inception-V3 (Szegedy et al., 2016), VGG-11/13/16/19 (Simonyan & Zisserman, 2015), ResNet-34/50/101/152 (He et al., 2016), EfficientNet-B0 (Tan & Le, 2019), MobileNet-V2 (Sandler et al., 2018), DLA (Yu et al., 2018), ResNeXt-29 (2x64d/4x64d/32x4d/8x64d) (Xie et al., 2017), SeNet-18 (Hu et al., 2018), and ConvMixer (Trockman & Kolter, 2023). For each architecture, we vary the training procedure to obtain 15 models across four adversarial training methods: PGD (Madry et al., 2018), TRADES (Zhang et al., 2019), PGD-SCORE, and TRADES-SCORE (Pang et al., 2022).

We train all models under both ℓ_∞ and ℓ_2 threat models with the following steps:

1. We use PGD adversarial training to train eight models with batch size $\in \{128, 512\}$, a learning rate $\in \{0.1, 0.05\}$, and weight decay $\in \{10^{-4}, 10^{-5}\}$. We also save the overall best hyperparameter choice. For the ℓ_2 threat model, we fix the learning rate to 0.1 since we observe that with ℓ_∞ , 0.1 is strictly better than 0.05.
2. Using the best hyperparameter choice, we train one model with PGD-SCORE, three with TRADES, and three with TRADES-SCORE. For TRADES and TRADES-SCORE, we take their β parameter from 0.1, 0.3, 1.0.

After training, we observe that some locally trained models exhibit inferior accuracy and/or robustness that is abnormally lower than others. The influence of inferior models on the correlation analysis is discussed in Appendix F. Finally, we filter out all models with an overall performance (accuracy + robustness) below 110. This threshold is determined to exclude only those evidently inferior models so that the size of model zoo (557 after filtering) is still large enough to ensure the generality and comprehensiveness of the conclusions drawn on it.

D. Additional Results

D.1. OOD Performance and Ranking

Table 4. Performance, evaluated with OODRobustBench, of state-of-the-art models trained on CIFAR10 to be robust to ℓ_2 attacks. Top 3 results under each metric are highlighted by **bold** and/or underscore. The column “OOD” gives the overall OOD robustness which is the mean of the robustness to OOD_d and OOD_t .

Method	Model	Accuracy		Robustness				Ranking (Rob.)	
		ID	OOD_d	ID	OOD_d	OOD_t	OOD	ID	OOD
BDM (Wang et al., 2023)	WRN-70-16	95.54	80.04	84.97	60.83	36.65	48.74	1	1
BDM (Wang et al., 2023)	WRN-28-10	95.16	79.28	83.69	59.39	35.04	47.21	2	2
FDA (Rebuffi et al., 2021)	WRN-70-16 (extra)	95.74	79.90	82.36	57.94	31.71	44.82	3	4
Uncovering (Gowal et al., 2021a)	WRN-70-16 (extra)	94.74	78.78	80.56	56.18	30.48	43.33	4	6
FDA (Rebuffi et al., 2021)	WRN-70-16 (DDPM)	92.41	75.95	80.42	56.82	<u>34.58</u>	<u>45.70</u>	5	3
RATIO (Augustin et al., 2020)	WRN-34-10 (extra)	93.97	77.40	78.81	54.71	31.62	43.16	6	7
FDA (Rebuffi et al., 2021)	WRN-28-10	91.79	75.26	78.79	55.63	33.32	44.48	7	5
PORT (Sehwag et al., 2022)	WRN-34-10	90.93	74.00	77.29	54.33	29.44	41.88	8	8
RATIO (Augustin et al., 2020)	WRN-34-10	92.23	76.43	76.27	52.83	29.25	41.04	9	11
HATE (Rade & Moosavi-Dezfooli, 2022)	PreActResNet-18	90.57	73.55	76.14	53.35	29.69	41.52	10	9
FDA (Rebuffi et al., 2021)	RreActResNet-18	90.33	72.96	75.87	52.21	30.06	41.14	11	10
Uncovering (Gowal et al., 2021a)	WRN-70-16	90.89	74.71	74.51	52.20	25.76	38.98	12	15
PORT (Sehwag et al., 2022)	ResNet-18	89.76	72.31	74.42	51.76	26.68	39.22	13	13
AWP (Wu et al., 2020)	WRN-34-10	88.51	71.23	73.66	51.53	27.50	39.52	14	12
RATIO (Augustin et al., 2020)	ResNet-50	91.07	74.24	72.99	49.32	28.72	39.02	15	14
PGD10 (Engstrom et al., 2019)	ResNet-50	90.83	73.85	69.25	46.65	17.71	32.18	16	16
Overfitting (Rice et al., 2020)	PreActResNet-18	88.67	71.27	67.69	44.76	18.58	31.67	17	17
DDN (Rony et al., 2019)	WRN-38-10	89.04	71.77	66.46	44.54	18.31	31.42	18	18
MMA (Ding et al., 2020)	WRN-28-4	88.00	72.32	66.09	43.79	16.52	30.15	19	20

Table 5. Performance, evaluated with OODRobustBench, of state-of-the-art models trained on ImageNet to be robust to ℓ_∞ attacks. Top 3 results under each metric are highlighted by **bold** and/or underscore. The column “OOD” gives the overall OOD robustness which is the mean of the robustness to OOD_d and OOD_t .

Method	Model	Accuracy		Robustness				Ranking (Rob.)	
		ID	OOD_d	ID	OOD_d	OOD_t	OOD	ID	OOD
Comprehensive (Liu et al., 2023)	Swin-L	78.92	45.84	59.82	23.59	29.88	26.74	1	1
Comprehensive (Liu et al., 2023)	ConvNeXt-L	78.02	44.74	58.76	23.35	30.10	26.72	2	2
Revisiting (Singh et al., 2023)	ConvNeXt-L-ConvStem	77.00	44.05	57.82	23.09	27.98	25.53	3	3
Comprehensive (Liu et al., 2023)	Swin-B	76.16	42.58	56.26	21.45	27.02	24.24	4	7
Revisiting (Singh et al., 2023)	ConvNeXt-B-ConvStem	75.88	42.29	56.24	21.77	27.89	24.83	5	5
Comprehensive (Liu et al., 2023)	ConvNeXt-B	76.70	43.06	56.02	21.74	26.97	24.36	6	6
Revisiting (Singh et al., 2023)	ViT-B-ConvStem	76.30	<u>44.67</u>	54.90	21.76	<u>28.98</u>	25.37	7	4
Revisiting (Singh et al., 2023)	ConvNeXt-S-ConvStem	74.08	39.55	52.66	19.35	26.87	23.11	8	9
Revisiting (Singh et al., 2023)	ConvNeXt-B	75.08	40.68	52.44	20.09	26.06	23.07	9	10
Comprehensive (Liu et al., 2023)	Swin-S	75.20	40.84	52.10	19.67	24.73	22.20	10	12
Comprehensive (Liu et al., 2023)	ConvNeXt-S	75.64	40.91	51.66	19.40	25.00	22.20	11	11
Revisiting (Singh et al., 2023)	ConvNeXt-T-ConvStem	72.70	38.15	49.46	17.97	25.32	21.65	12	14
Revisiting (Singh et al., 2023)	ViT-S-ConvStem	72.58	39.24	48.46	17.83	25.43	21.63	13	15
Revisiting (Singh et al., 2023)	ViT-B	72.98	42.38	48.34	20.43	26.26	23.34	14	8
Light (Debenedetti et al., 2023)	XCiT-L12	73.78	38.10	47.88	15.84	23.22	19.53	15	18
Revisiting (Singh et al., 2023)	ViT-M	71.78	39.88	47.34	18.95	25.25	22.10	16	13
Revisiting (Singh et al., 2023)	ConvNeXt-T	71.88	37.70	46.98	17.13	21.36	19.25	17	19
Easy (Mao et al., 2022)	Swin-B	74.14	38.45	46.54	15.36	22.19	18.78	18	20
Comprehensive (Liu et al., 2023)	ViT-B	72.84	39.88	45.90	18.01	22.95	20.48	19	16
Light (Debenedetti et al., 2023)	XCiT-M12	74.04	37.00	45.76	14.73	22.82	18.77	20	21

D.2. Performance Degradation Distribution

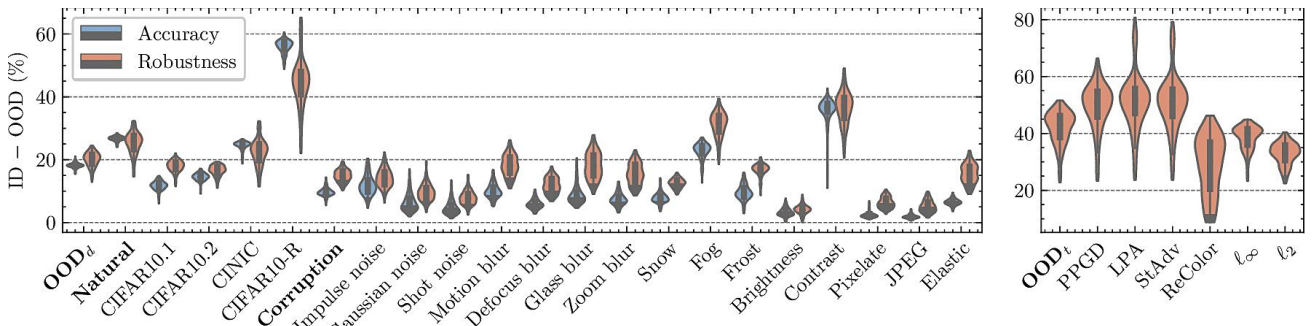


Figure 9. Degradation of accuracy and robustness under various distribution shifts for CIFAR10 l_2 .

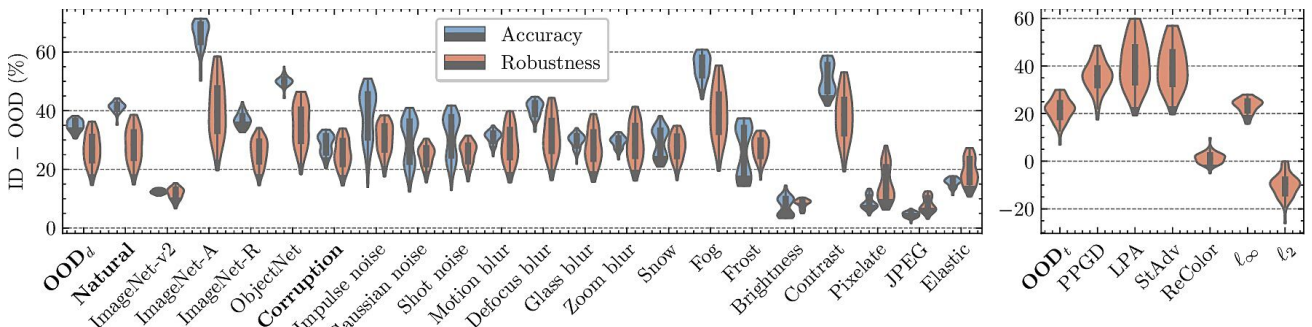


Figure 10. Degradation of accuracy and robustness under various distribution shifts for ImageNet l_∞ .

D.3. Correlation Between ID and OOD Performance under Dataset Shifts

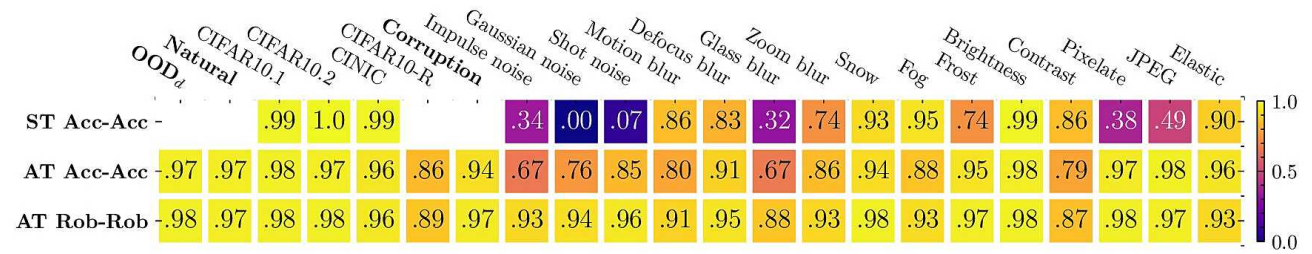


Figure 11. R^2 of regression between ID and OOD performance for Standardly-Trained (ST) and Adversarially-Trained (AT) models under dataset shifts for CIFAR10 l_2 . Higher R^2 implies stronger linear correlation. The result of ST is copied from (Miller et al., 2021).

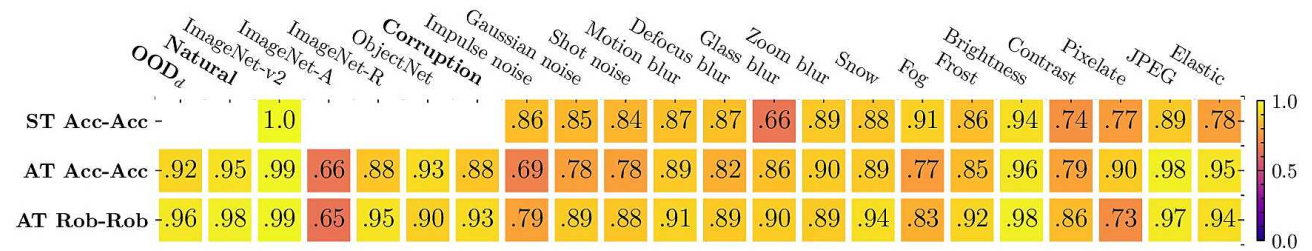


Figure 12. R^2 of regression between ID and OOD performance for Standardly-Trained (ST) and Adversarially-Trained (AT) models under dataset shifts for ImageNet l_∞ . Higher R^2 implies stronger linear correlation. The result of ST is copied from (Miller et al., 2021).

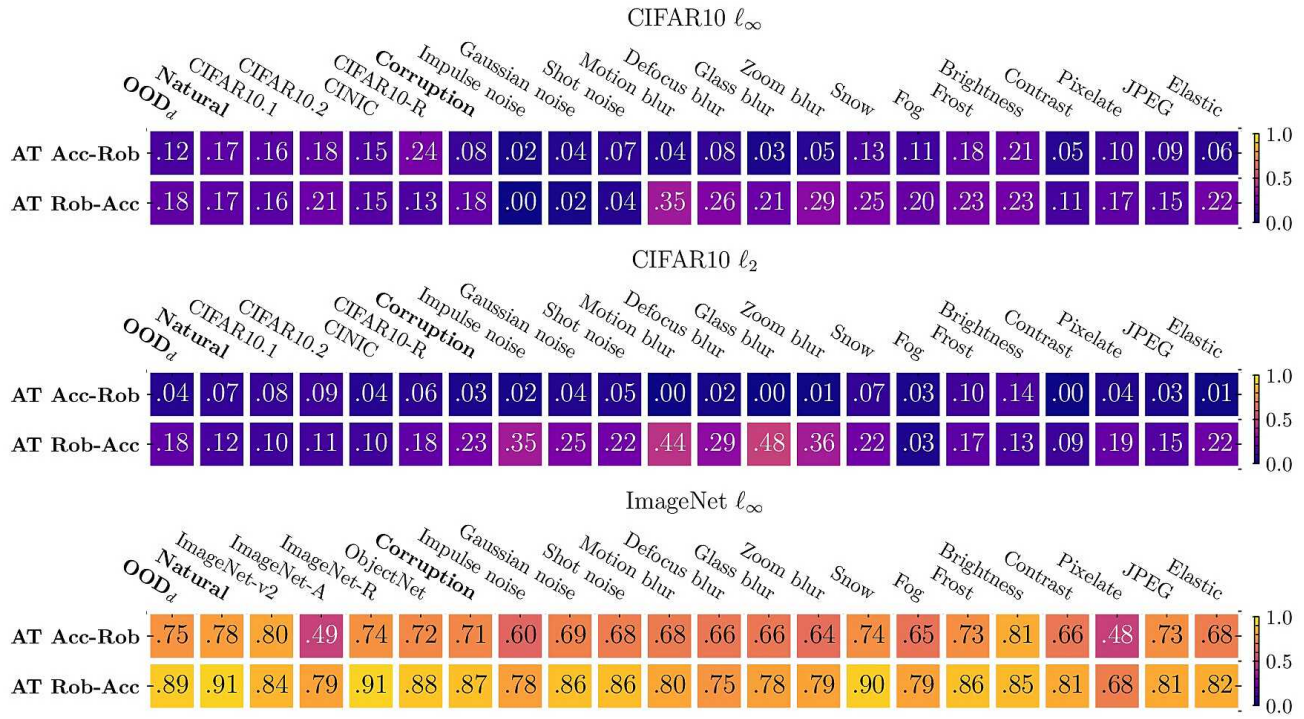


Figure 13. R^2 of regression between ID and OOD performance for Adversarially-Trained (AT) models under various dataset shifts. "Acc-Rob" denotes the linear model between ID accuracy (x) and OOD robustness (y) and "Rob-Acc" for ID robustness (x) and OOD accuracy (y).

D.4. Unsupervised OOD Robustness Prediction

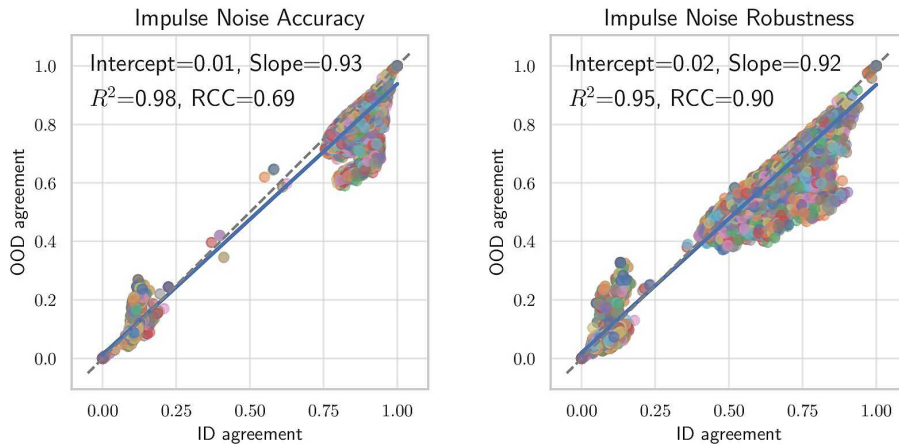


Figure 14. Correlation between ID and OOD prediction agreement on adversarial examples for CIFAR10 l_∞ AT models.

D.5. Predicted Upper Limit of OOD Accuracy and Robustness

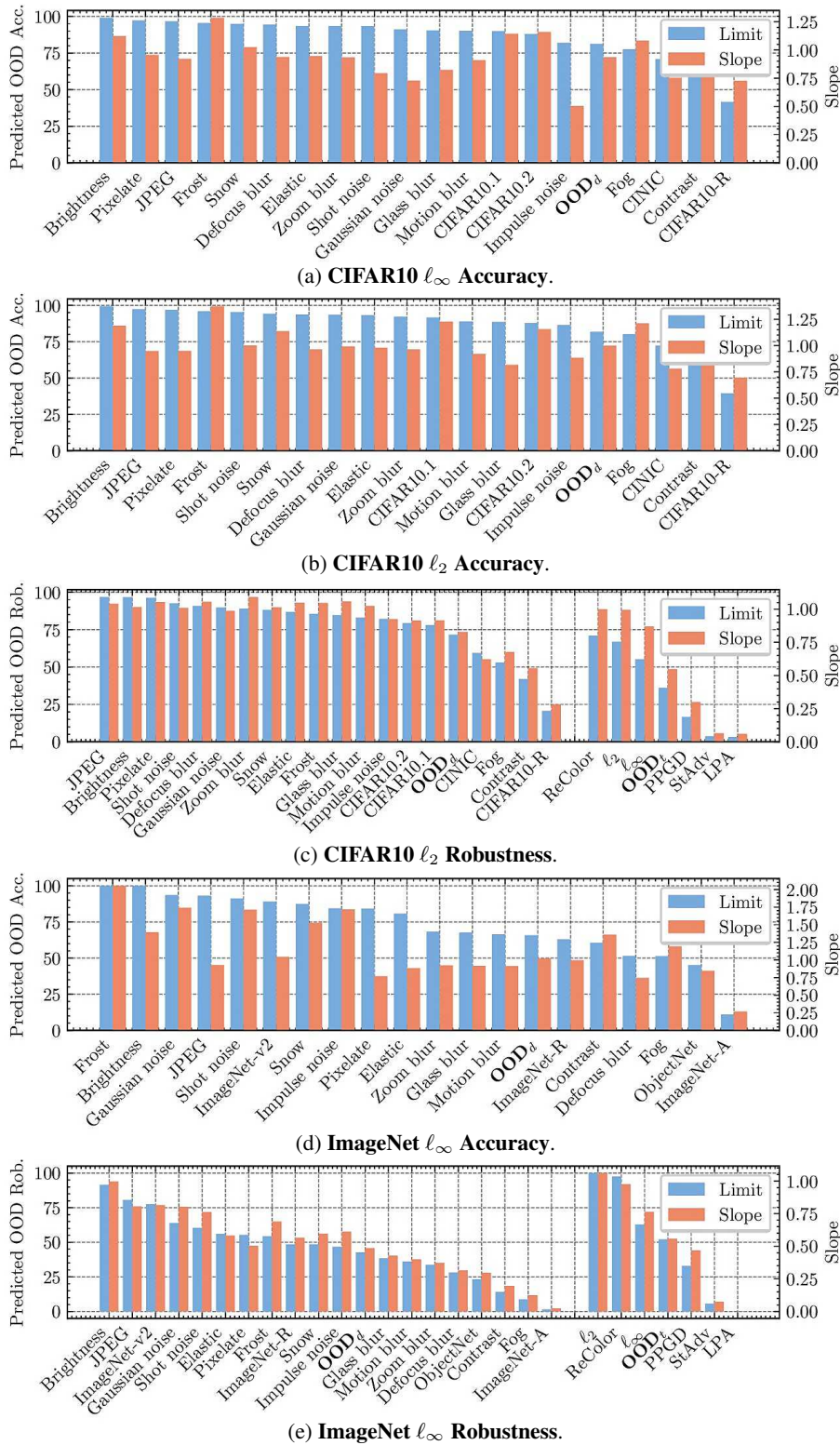


Figure 15. The estimated upper limit of OOD performance and the conversion rate, a.k.a. slope, to OOD performance from ID performance under various distribution shifts.

E. Catastrophic Degradation of Robustness

We observe this issue on only one implementation, using WideResNet28-10 with extra synthetic data (model id: *Rade2021Helper_ddpm* on RobustBench), from Rade & Moosavi-Dezfooli (2022) for CIFAR10 ℓ_∞ . There are three other implementations of this method on RobustBench. None of them, including the one using ResNet18 with extra synthetic data, is observed to suffer from this issue. It seems that catastrophic degradation in this case is specific to the implementation or training dynamics.

On the other hand, catastrophic degradation consistently happens on the models trained with AutoAugment or IDBH but not other tested data augmentations. It suggests the possibility that a certain image transformation operation exclusively used by AutoAugment and IDBH cause this issue. Besides, catastrophic degradation also consistently happens on the models trained using the receipt of (Debenedetti et al., 2023) under Gaussian and shot noise shifts. However, it employs a wide range of training techniques, so further experiments are required to identify the specific cause.

F. How Inferior Models Affect the Correlation Analysis

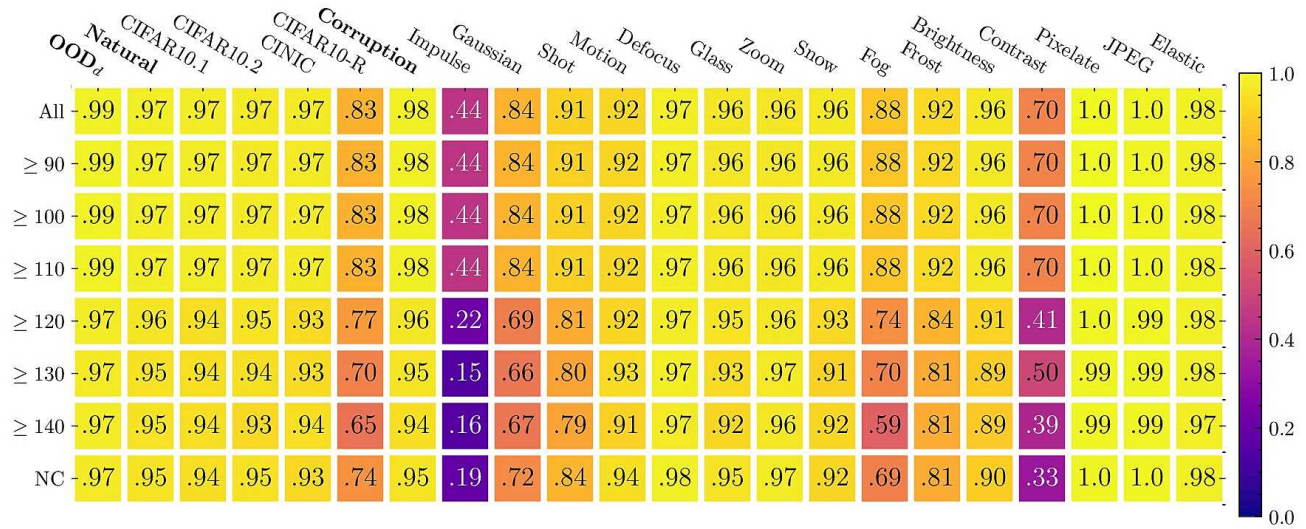
This section studies the influence of the construction of model zoo on the result of correlation. We use the overall performance (accuracy + robustness) to filter out inferior models. As we increasing the threshold of overall performance for filtering, the average overall performance of the model zoo increases, the number of included models decreases and the weight of the models from other published sources on the regression grows up. Our locally trained models are normally inferior to the public models regarding the performance since the latter employs better optimized and more effective training methods and settings. The training methods and settings of public models are also much more diverse.

The correlation for particular shifts varies considerably as more inferior models removed. R^2 declines considerably under CIFAR10-R, noise, fog, glass blur, frost and contrast for both Acc-Acc and Rob-Rob on CIFAR10 ℓ_∞ (Figure 16) and ℓ_2 (Figure 17). A similar trend is also observed for threat shifts, ReColor and different p -norm for CIFAR10 ℓ_∞ as shown in Figure 18. It suggests that the weak correlation under these shifts mainly results from those high-performance public models, and is likely related to the fact that these models include much diverse training methods and settings. For example, all observed catastrophic degradation under the noise shifts occur in the public models. Note that the locally trained models have a large diversity in model architectures particularly within the family of CNNs, but it seems that this architectural diversity does not effect the correlation as much as other factors.

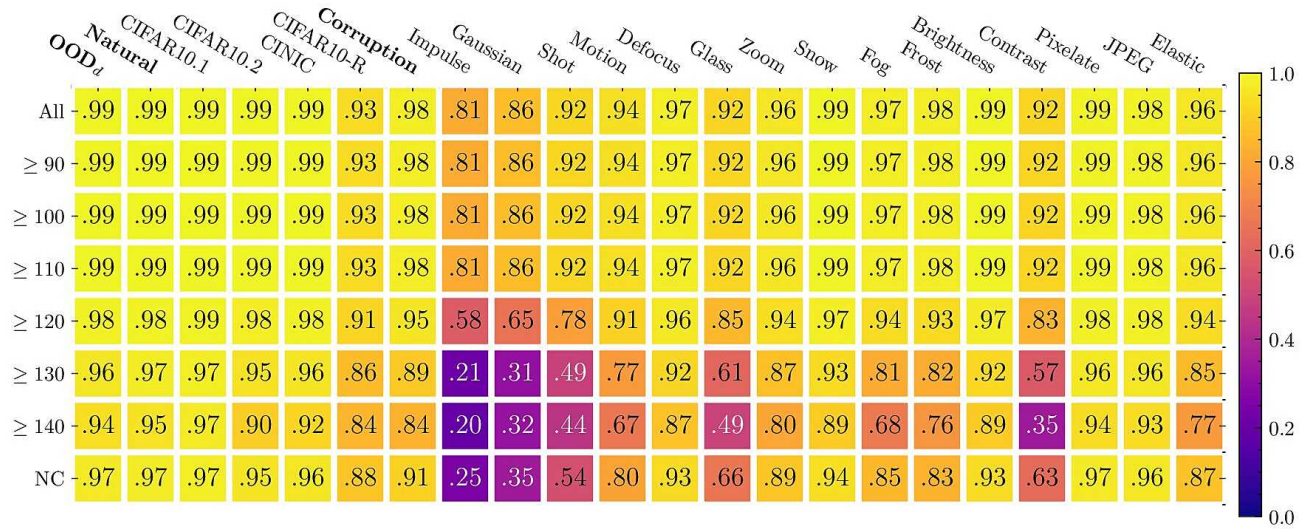
In contrast, correlation is improved for most threat shifts for CIFAR10 ℓ_2 as shown in Figure 18. As shown in Figure 29, the locally trained (inferior) models and the public (high-performance) models have divergent linear trends (most evident in the plot of PPGD). That’s why removing models from either group will enhance the correlation. Note that such divergence is not evident in the figures of CIFAR10 ℓ_∞ (Figure 28) and ImageNet ℓ_∞ (Figure 30).

F.1. No Evident Correlation when ID and OOD Metrics Misalign

Inferior models also cause OOD robustness to not consistently increase with the ID accuracy, i.e., the poor correlation between ID accuracy (robustness) and OOD robustness (accuracy) because they have high accuracy yet poor robustness. These models are mainly produced by some of our custom training receipts and take a considerable proportion of our CIFAR-10 model zoo, whereas the model zoo of ImageNet is dominated by ones from public sources.

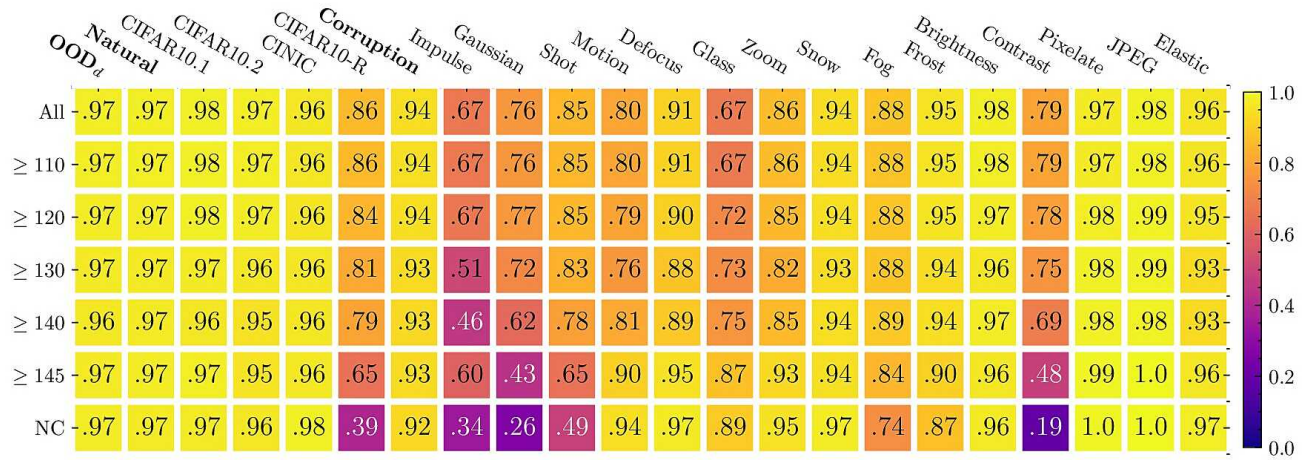


(a) R^2 of Acc-Acc.



(b) R^2 of Rob-Rob.

Figure 16. The change of R^2 under various dataset shifts as the models with lower overall performance are removed from regression for CIFAR10 ℓ_∞ . Each row, with the filtering threshold labeled at the lead, corresponds to a new filtered model zoo and the regression conducted it. "NC" refers to No Custom models, so all models are retrieved from either RobustBench or other published works.



(a) R^2 of Acc-Acc.



(b) R^2 of Rob-Rob.

Figure 17. The change of R^2 under various dataset shifts as the models with lower overall performance are removed from regression for CIFAR10 ℓ_2 . Each row, with the filtering threshold labeled at the lead, corresponds to a new filtered model zoo and the regression conducted it. "NC" refers to No Custom models, so all models are retrieved from either RobustBench or other published works.

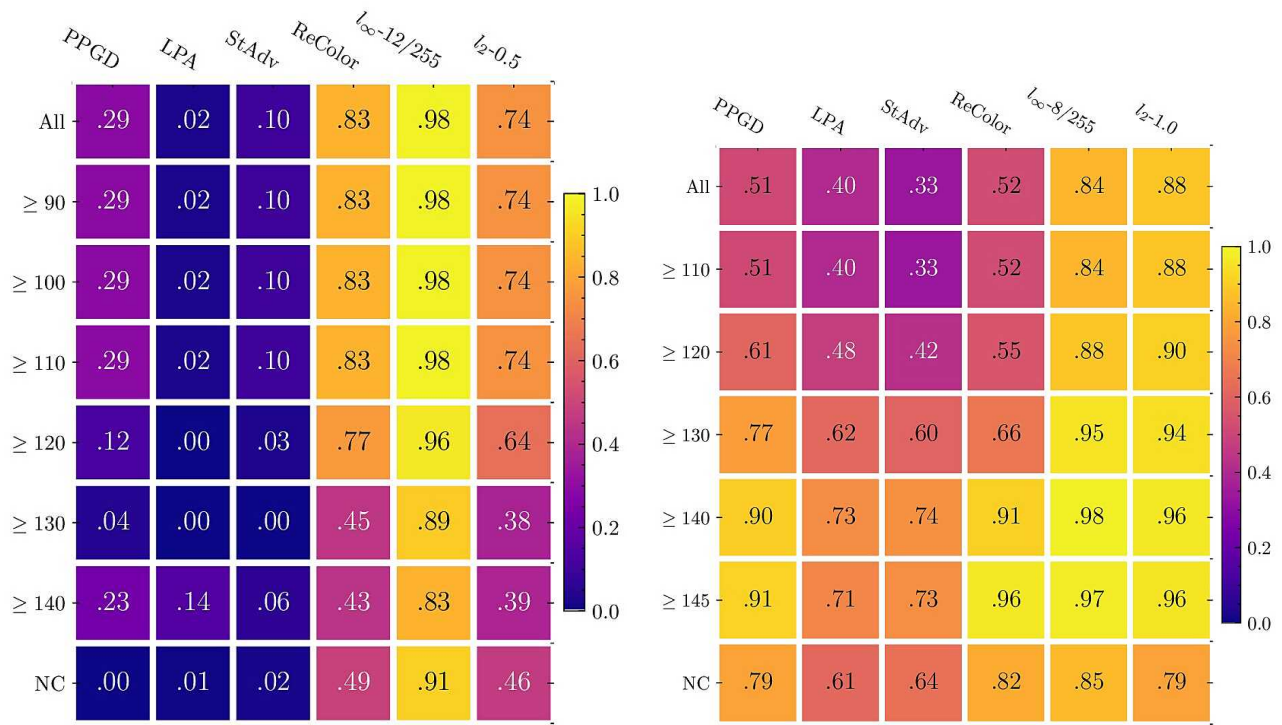


Figure 18. The change of R^2 under various threat shifts as the models with lower overall performance are removed from regression. Each row, with the filtering threshold labeled at the lead, corresponds to a new filtered model zoo and the regression conducted on it. "NC" refers to No Custom models, so all models are retrieved from either RobustBench or other published works.

G. Methods for Improving OOD Adversarial Robustness

All models used in this analysis are retrieved from RobustBench or other published works to ensure they are well-trained by the techniques to be examined. The specific experiment setting for each model can be found in its original paper.

Table 6. The effect of training with extra data on the OOD generalization of accuracy and robustness.

Dataset	Threat Model	Training	Model Architecture	Extra Data	ID		OOD _d				OOD _t	
					Acc.	Rob.	Acc.	Rob.	EAcc.	ERob.	Rob.	ERob.
CIFAR10	Linf	(Gowal et al., 2021a)	WideResNet70-16	-	85.29	57.24	66.98	35.90	-0.56	0.30	29.39	-2.18
				Synthetic	88.74	66.24	70.68	42.76	-0.08	0.74	33.65	-2.13
				Real	91.10	66.03	73.24	42.58	0.26	0.71	34.00	-1.67

Table 7. The effect of data augmentation on the OOD generalization of accuracy and robustness. The results reported in Figure 7b are the mean of the results on ViT and WideResNets.

Dataset	Threat Model	Training	Model Architecture	Data Augmentation	ID		OOD _d				OOD _t	
					Acc.	Rob.	Acc.	Rob.	EAcc.	ERob.	Rob.	ERob.
CIFAR10	Linf	(Li & Spratling, 2023c)	ViT-B	RandomCrop	83.23	47.02	66.48	28.85	0.86	0.54	27.36	0.57
				Cutout	84.22	49.57	67.23	30.68	0.69	0.56	29.74	1.75
				CutMix	80.92	47.45	63.93	29.89	0.48	1.27	30.48	3.49
				TrivialAugment	80.33	46.61	64.59	29.56	1.69	1.54	30.40	3.80
				AutoAugment	82.75	48.11	65.89	29.78	0.73	0.69	30.90	3.60
				IDBH	86.92	51.55	70.51	32.08	1.45	0.54	30.59	1.68
				IDBH	86.92	51.55	70.51	32.08	1.45	0.54	30.59	1.68
			WideResNet34-10	RandomCrop	86.52	52.42	68.11	31.55	-0.58	-0.61	26.47	-2.84
				Cutout	86.77	53.31	68.40	31.03	-0.53	-1.76	27.00	-2.74
				CutMix	87.41	53.89	68.97	31.71	-0.55	-1.50	28.50	-1.50
				TrivialAugment	86.98	54.18	69.85	32.94	0.73	-0.47	28.62	-1.52
				AutoAugment	87.93	55.10	70.05	32.17	0.04	-1.90	29.06	-1.51
				IDBH	88.62	55.56	70.96	32.99	0.30	-1.41	28.58	-2.21
				IDBH	88.62	55.56	70.96	32.99	0.30	-1.41	28.58	-2.21

Table 8. The effect of model architecture on the OOD generalization of accuracy and robustness.

Dataset	Threat Model	Training	Model Architecture	Model Size (M)	ID		OOD _d				OOD _t	
					Acc.	Rob.	Acc.	Rob.	EAcc.	ERob.	Rob.	ERob.
ImageNet	ℓ_∞	(Liu et al., 2023)	ResNet152	60.19	70.92	43.62	34.43	14.13	-1.71	-1.26	17.23	-3.47
			ConvNeXt-B	88.59	76.70	56.02	43.06	21.74	1.03	0.33	26.97	-0.63
			ViT-B	86.57	72.84	45.90	39.88	18.01	1.78	1.51	22.95	0.98
			Swin-B	87.77	76.16	56.26	42.58	21.45	1.10	-0.07	27.02	-0.72

Table 9. The effect of model size on the OOD generalization of accuracy and robustness. The results reported in Figure 7d are averaged over three architectures at the corresponding relatively model size. For example, the result of "small" is averaged over WideResNet28-10, ResNet50 and ConvNeXt-S-ConvStem.

Dataset	Threat Model	Training	Model Architecture	Model Size	ID		OOD _d				OOD _t	
					Acc.	Rob.	Acc.	Rob.	EAcc.	ERob.	Rob.	ERob.
CIFAR10	ℓ_∞	(Rebuffi et al., 2021)	WideResNet28-10	36.48	87.33	60.88	69.35	38.54	-0.10	0.35	33.63	0.36
			WideResNet70-16	266.80	88.54	64.33	70.62	41.01	0.04	0.35	34.12	-0.76
			WideResNet106-16	415.48	88.50	64.82	70.65	41.43	0.11	0.42	33.90	-1.22
ImageNet	ℓ_∞	(Liu et al., 2023)	ResNet50	25.56	65.02	32.02	28.43	9.23	-1.68	-0.53	13.71	-0.52
			ResNet101	44.55	68.34	39.76	31.74	12.44	-1.76	-1.08	16.82	-1.72
			ResNet152	60.19	70.92	43.62	34.43	14.13	-1.71	-1.26	17.23	-3.47
ImageNet	ℓ_∞	(Singh et al., 2023)	ConvNeXt-S-ConvStem	50.26	74.08	52.66	39.55	19.35	0.19	-0.42	26.87	1.14
			ConvNeXt-B-ConvStem	88.75	75.88	56.24	42.29	21.77	1.10	0.26	27.89	0.16
			ConvNeXt-L-ConvStem	198.13	77.00	57.82	44.05	23.09	1.71	0.80	27.98	-0.63

Table 10. The effect of different adversarial training methods on the OOD generalization of accuracy and robustness.

Dataset	Threat	Training	ID		OOD _d				OOD _t	
			Acc.	Rob.	Acc.	Rob.	EAcc.	ERob.	Rob.	ERob.
CIFAR10	ℓ_∞	PGD (Li & Spratling, 2023c)	86.52	52.42	68.11	31.55	-0.58	-0.61	26.47	-2.84
		VR- ℓ_∞ (Dai et al., 2022)	72.72	49.92	56.12	31.84	0.34	1.47	34.70	6.55
		PGD (Rice et al., 2020)	85.34	53.52	66.46	32.07	-1.12	-0.88	27.89	-1.94
		HE (Pang et al., 2020)	85.14	53.84	66.96	32.45	-0.43	-0.72	46.20	16.22
		PGD (locally-trained)	80.44	38.98	62.40	22.18	-0.60	-0.39	21.77	-1.27
		MMA (Ding et al., 2020)	84.37	41.86	68.22	24.65	1.54	0.02	35.12	10.74
		PGD Goyal et al. (2021a)	91.10	66.03	73.24	42.58	0.26	0.71	34.00	-1.67
AS (Bai et al., 2023)	95.23	69.50	79.09	43.32	2.25	-1.03	46.71	9.41		

H. Plots of ID-OOD Correlation per Dataset Shift

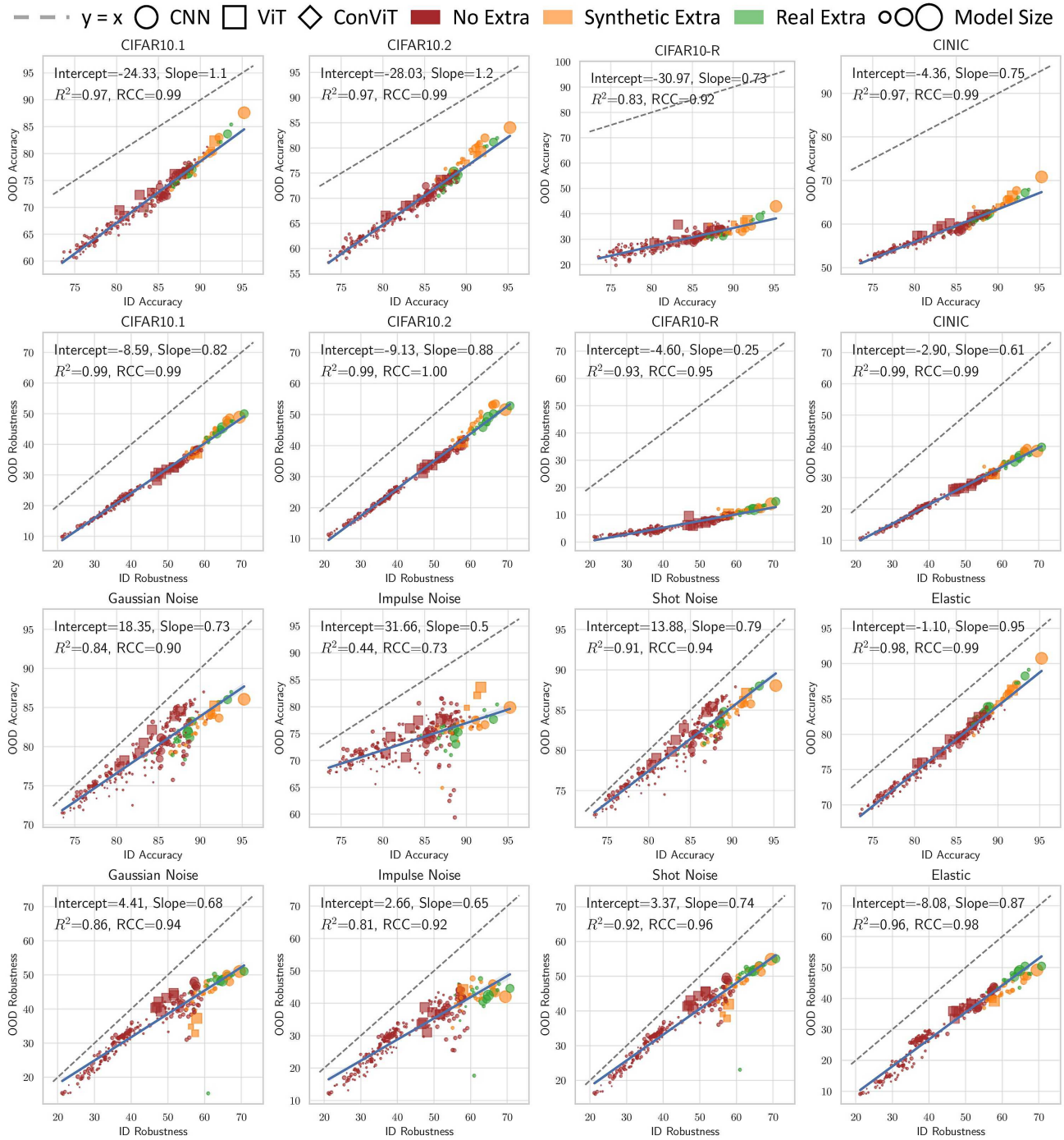


Figure 19. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for CIFAR10 ℓ_∞ AT models.

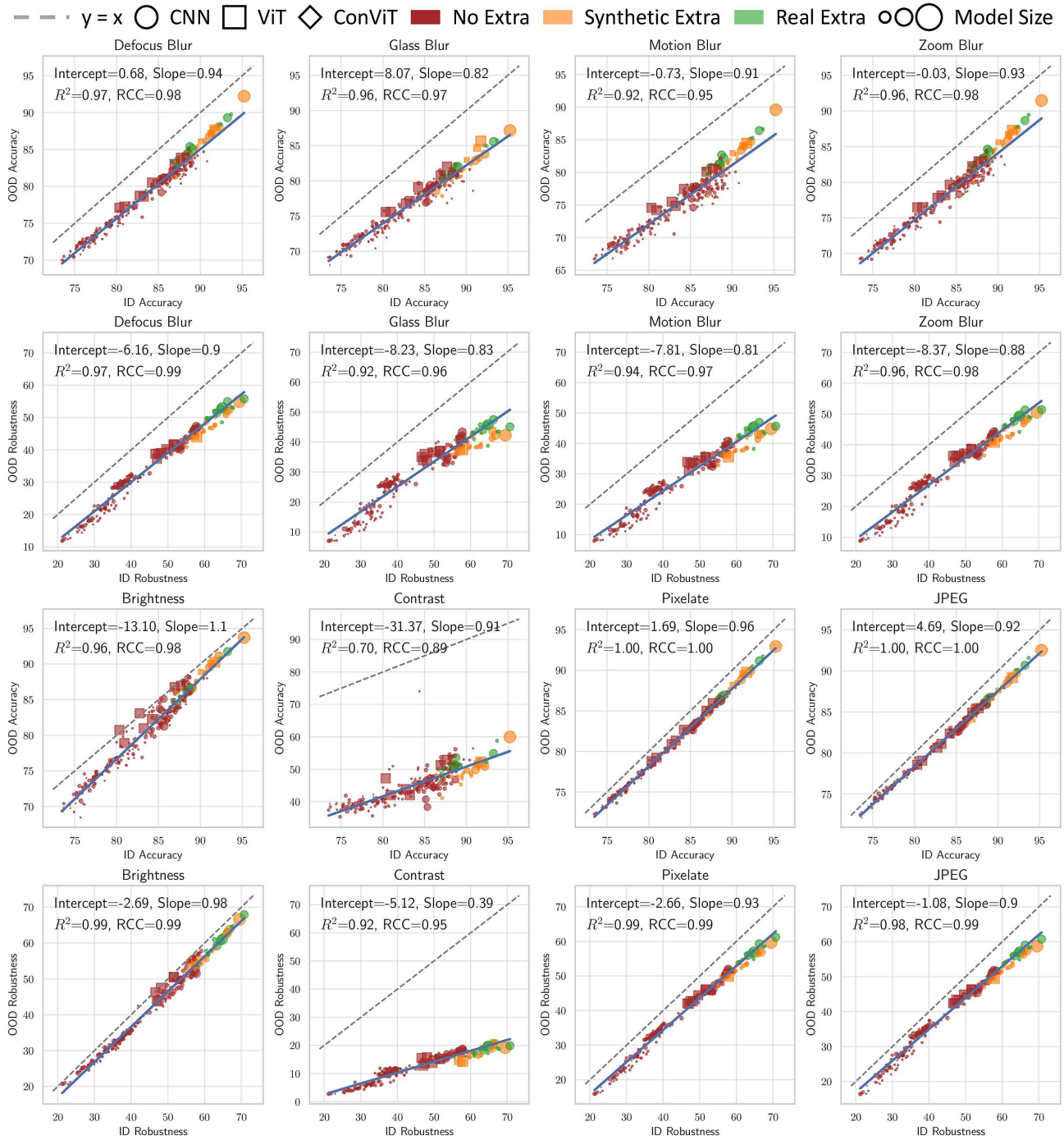


Figure 20. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for CIFAR10 ℓ_∞ AT models.

--- $y = x$ ○ CNN □ ViT ◇ ConViT ■ No Extra ■ Synthetic Extra ■ Real Extra ○○○○ Model Size

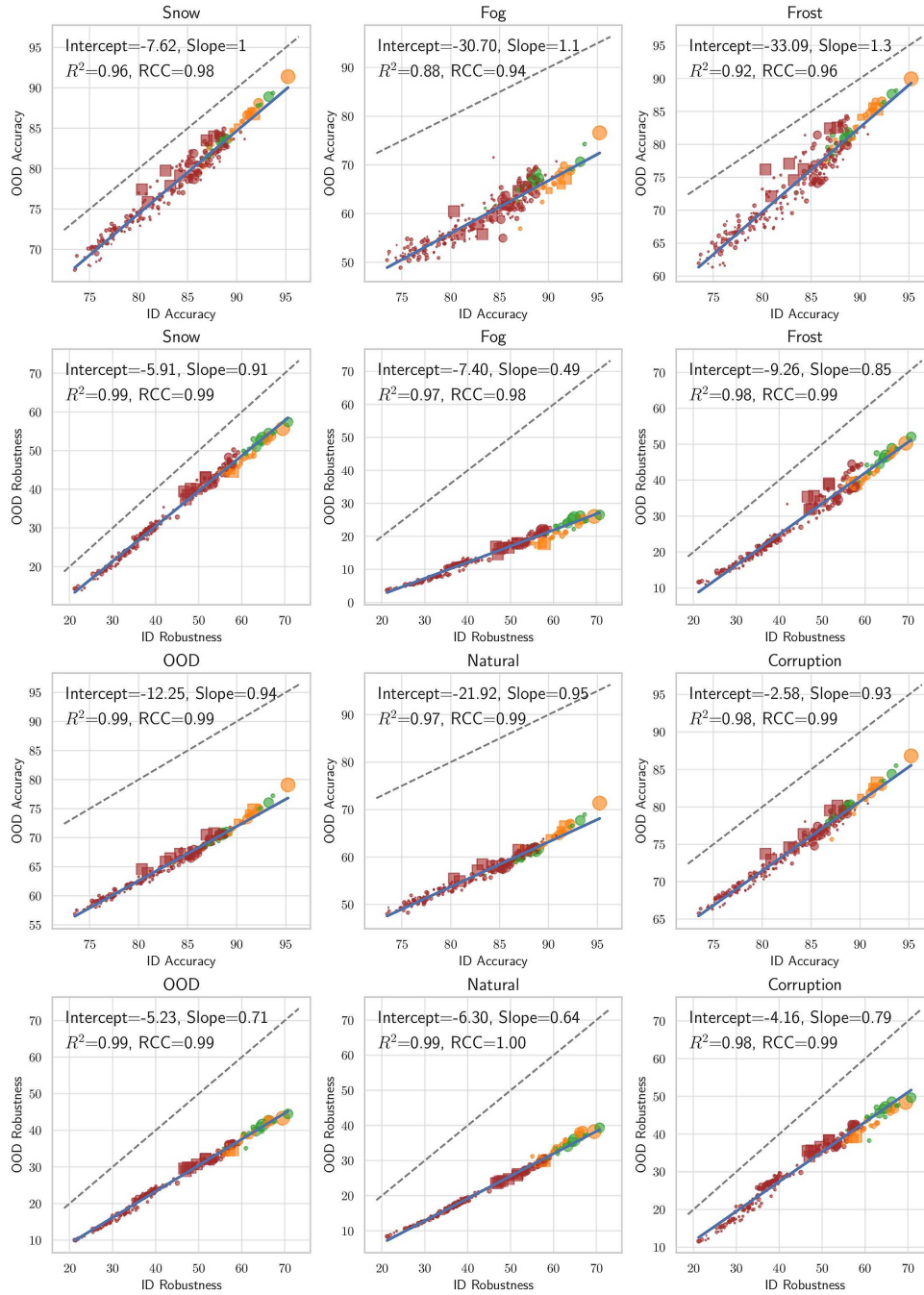


Figure 21. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for CIFAR10 ℓ_∞ AT models.

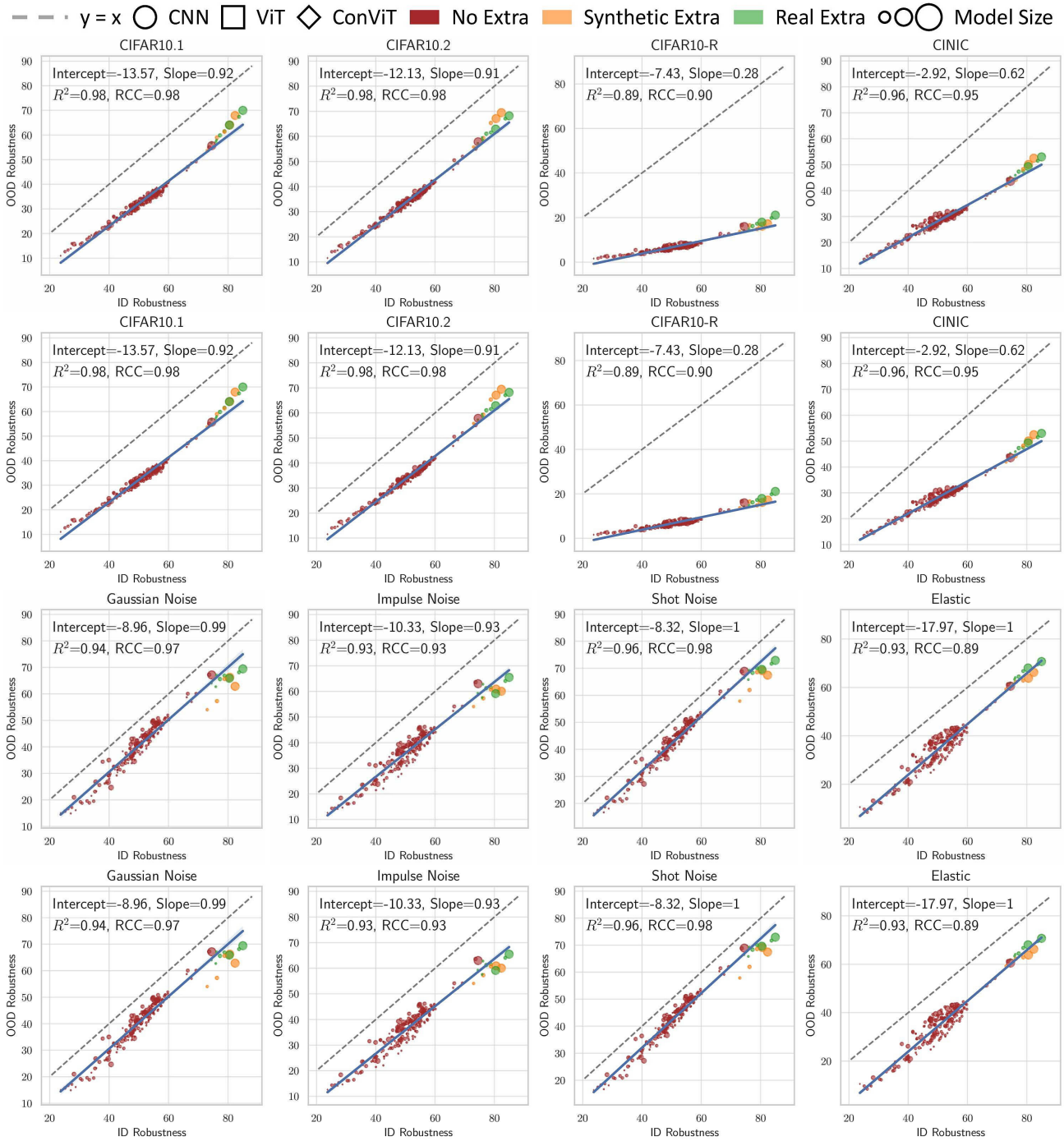


Figure 22. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for CIFAR10 ℓ_2 AT models.

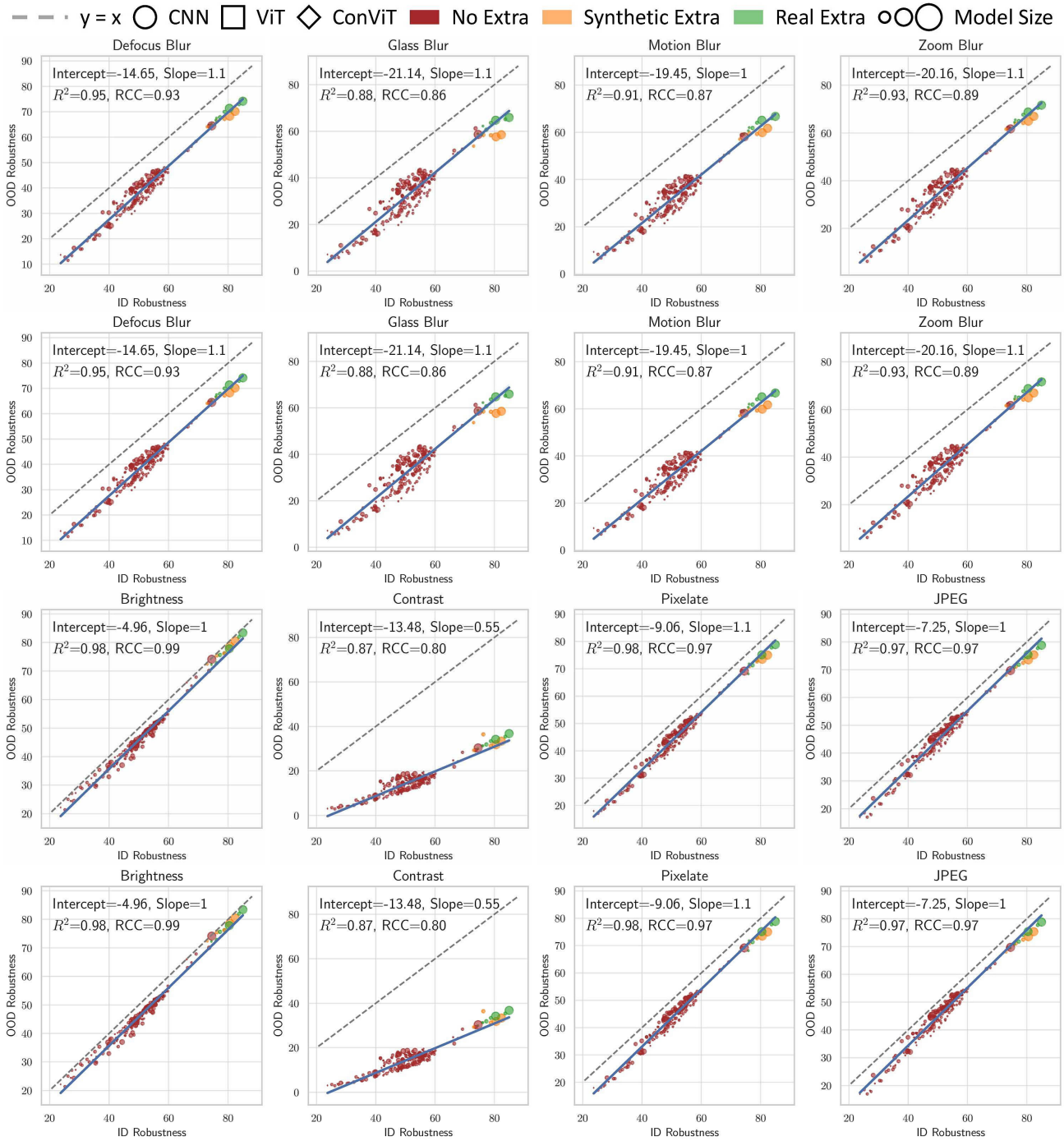


Figure 23. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for CIFAR10 l_2 AT models.

--- $y = x$ ○ CNN □ ViT ◇ ConViT ■ No Extra ■ Synthetic Extra ■ Real Extra ○○○○ Model Size

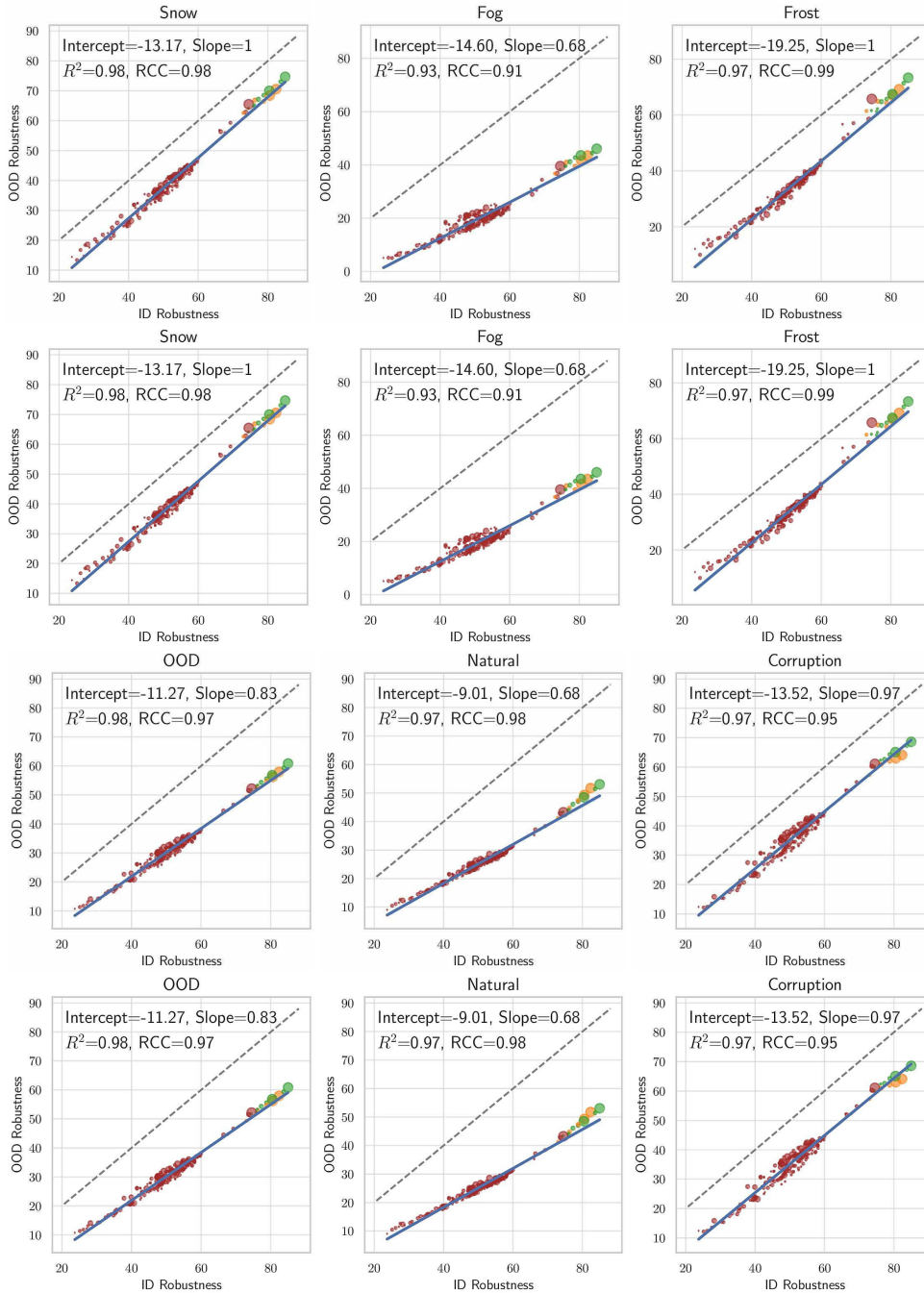


Figure 24. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for CIFAR10 ℓ_2 AT models.

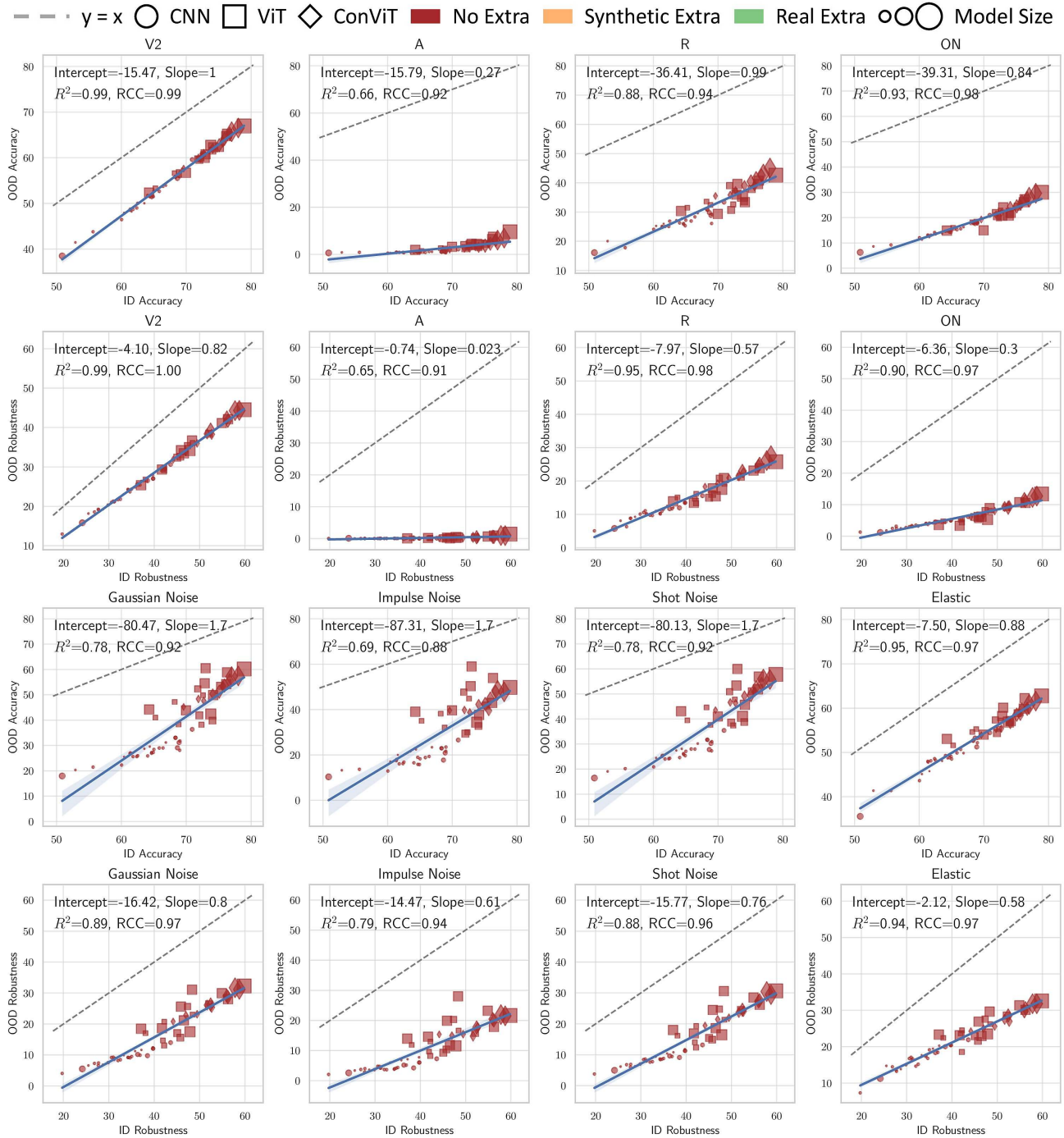


Figure 25. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for ImageNet ℓ_∞ AT models.

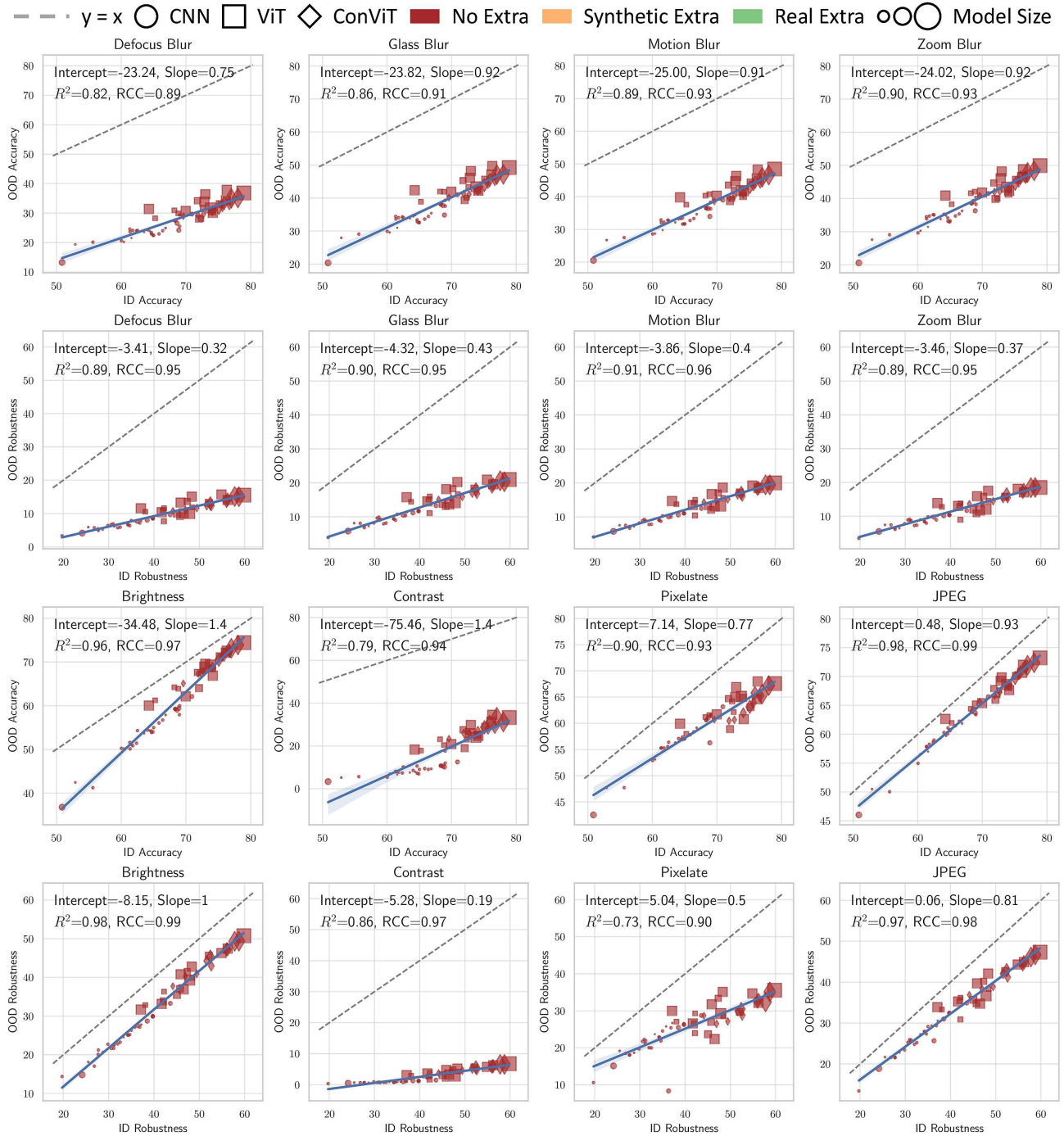


Figure 26. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for ImageNet ℓ_∞ AT models.

--- $y = x$ ○ CNN □ ViT ◇ ConViT ■ No Extra ■ Synthetic Extra ■ Real Extra ○○○ Model Size

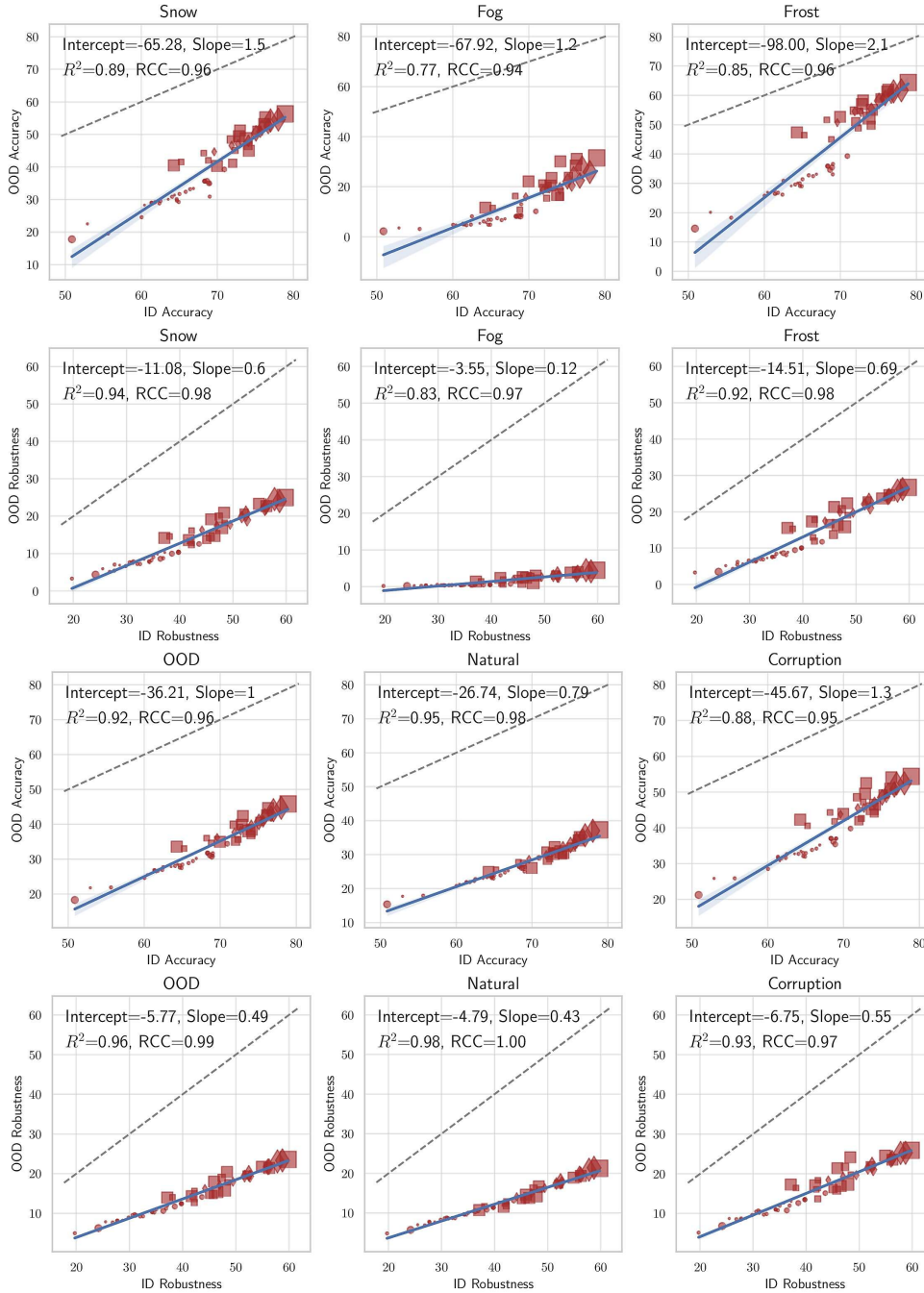


Figure 27. Correlation between ID accuracy and OOD accuracy (odd rows); ID robustness and OOD robustness (even rows) for ImageNet ℓ_∞ AT models.

I. Plots of ID-OOD Correlation per Threat Shift

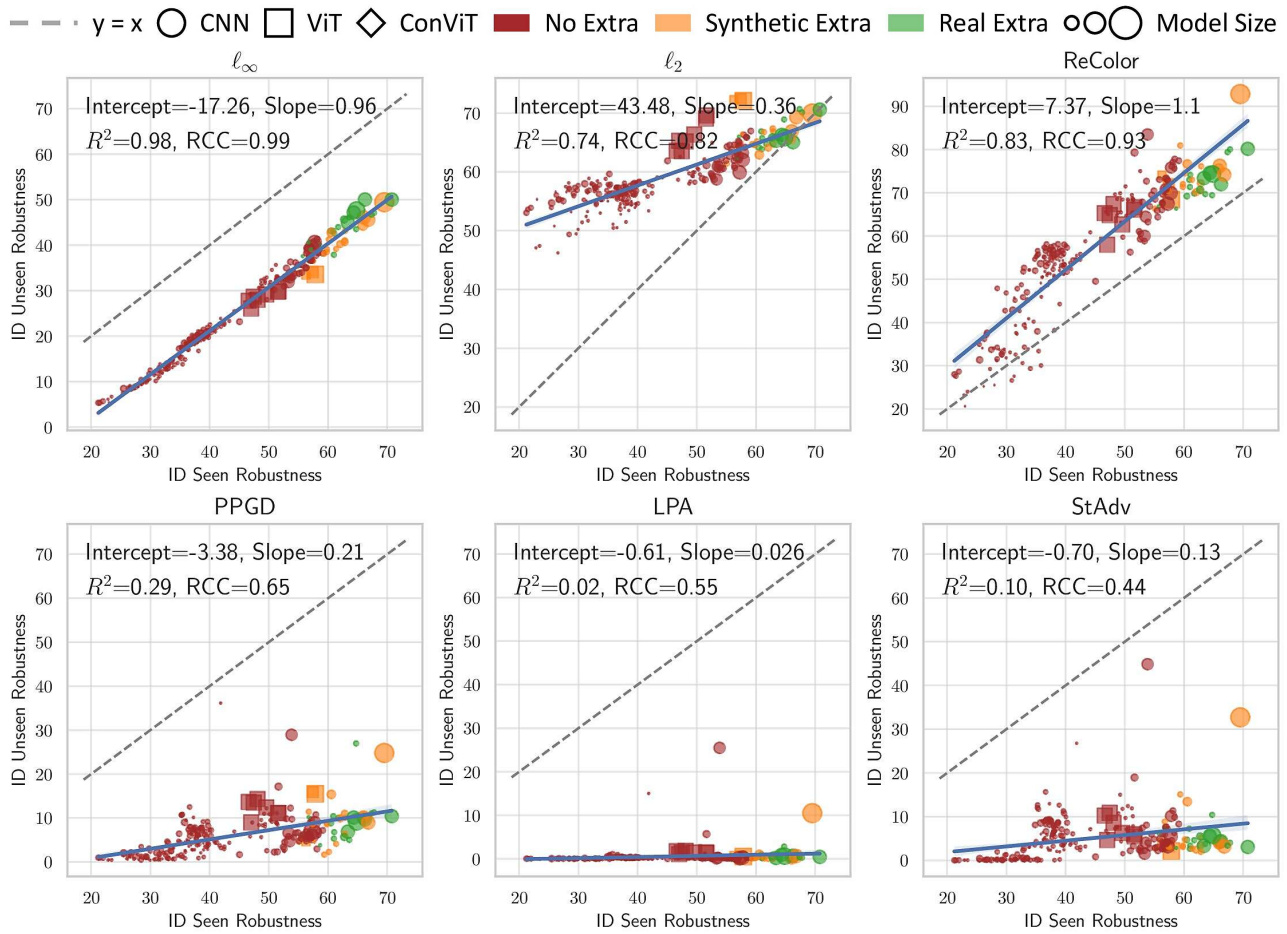


Figure 28. Correlation between seen and unforeseen robustness on ID data for CIFAR10 l_∞ AT models.

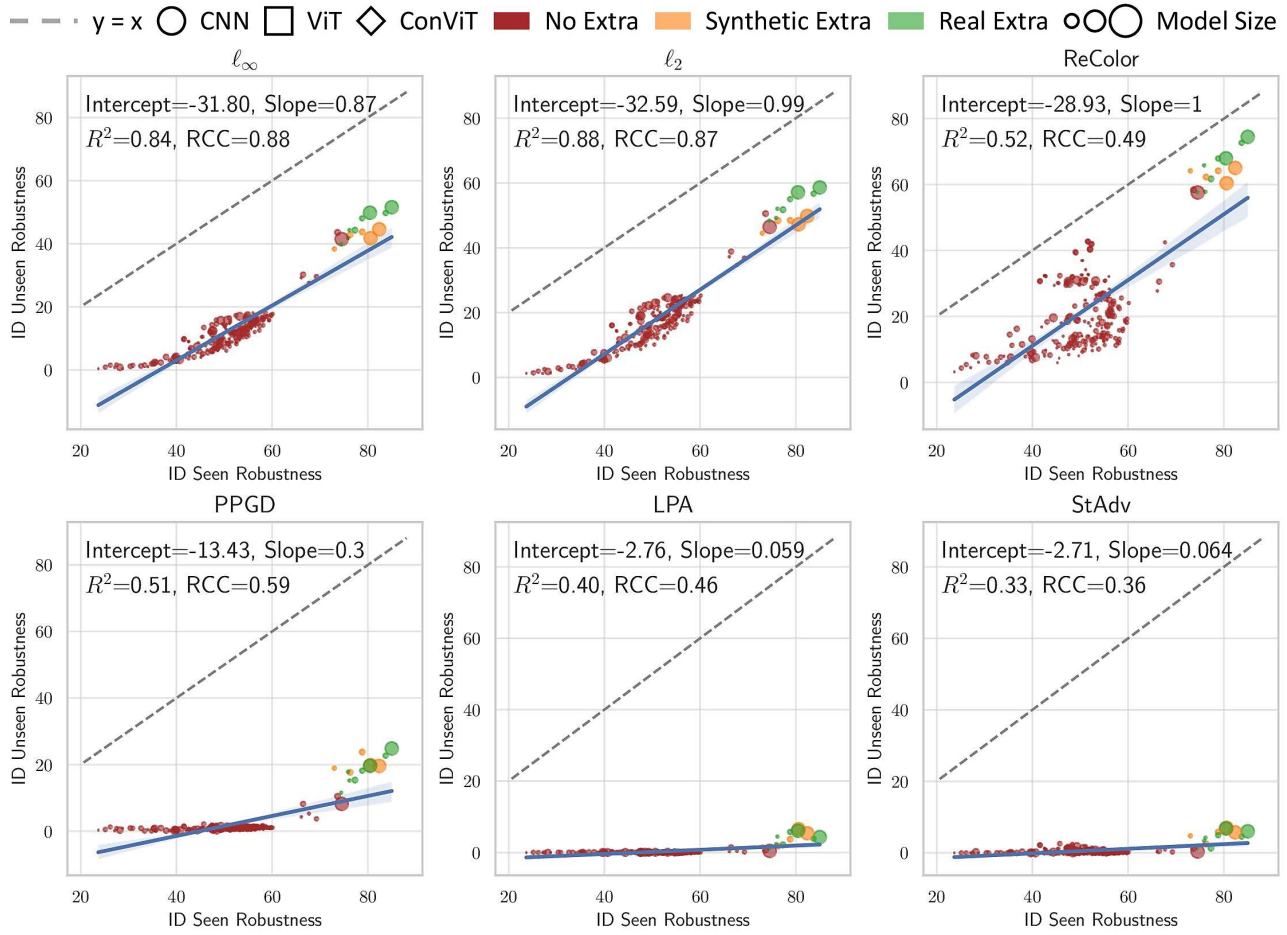


Figure 29. Correlation between seen and unforesen robustness on ID data for CIFAR10 l_2 AT models.

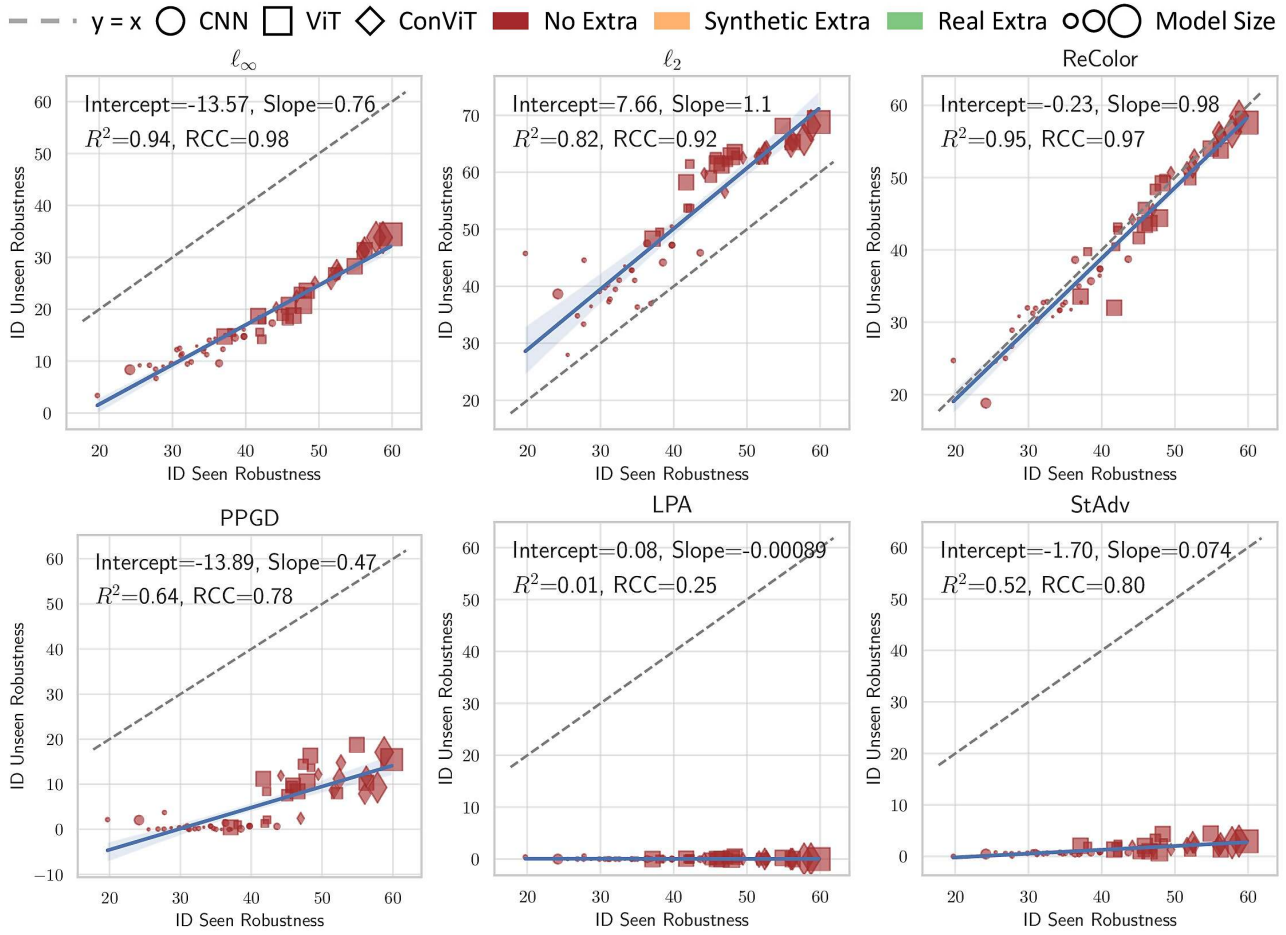


Figure 30. Correlation between seen and unforeseen robustness on ID data for ImageNet ℓ_∞ AT models.