

EVALUATING MODEL ROBUSTNESS AGAINST UNFORESEEN ADVERSARIAL ATTACKS

Anonymous authors

Paper under double-blind review

ABSTRACT

When considering real-world adversarial settings, defenders are unlikely to have access to the full range of deployment-time adversaries, and adversaries are likely to use realistic adversarial distortions that will not be limited to small L_p -constrained perturbations. To narrow in on this discrepancy between research and reality we introduce ImageNet-UA, a new benchmark for evaluating model robustness against a wide range of *unforeseen adversaries*. We make use of our benchmark to identify holes in current popular adversarial defense techniques, highlighting a rich space of techniques which can improve unforeseen robustness. We hope the greater variety and realism of ImageNet-UA will make it a useful tool for those working on real-world worst-case robustness, enabling development of more robust defenses which can generalize beyond attacks seen during training.

1 INTRODUCTION

Neural networks perform well on a variety of tasks, yet can be consistently fooled by minor adversarial distortions (Szegedy et al., 2013; Goodfellow et al., 2014b). This has led to an extensive area of research around the “ L_p -bounded adversary” that adds imperceptible distortions to model inputs to cause misclassification. However, this classic threat model may fail to fully capture many real-world concerns regarding worst-case robustness (Gilmer et al., 2018). Firstly, real-world worst-case distributions are likely to be varied, and are unlikely to be constrained to the L_p ball. Secondly, developers will not have access to the worst-case inputs to which their systems will be exposed to. For example, online advertisers use perturbed pixels in ads to defeat ad blockers trained only on the previous generation of ads in an ever-escalating arms race (Tramèr et al., 2018). Furthermore, although research has shown that adversarial training can lead to overfitting, wherein robustness against one particular adversary does not generalize (Dai et al., 2022; Yu et al., 2021; Stutz et al., 2020; Tramer & Boneh, 2019), the existing literature is still focuses on defenses that train against the test-time attacks. Although such distribution shifts have been studied in the average-case common corruption setting (Hendrycks & Dietterich, 2018), when considering worst-case inputs, the research community is lacking a unified benchmark for testing how defences generalise.

We address the limitations of current adversarial robustness evaluations by providing a repository of nineteen gradient-based attacks, eight of which are used to create ImageNet-UA—a benchmark for evaluating the *unforeseen robustness* of models on the popular ImageNet dataset (Deng et al., 2009). Defenses achieving high Unforeseen Adversarial Accuracy (UA2) on ImageNet-UA generalize to a diverse set of adversaries not seen at train time, demonstrating robustness to a much more realistic threat model than the L_p adversaries which are a focus of the literature.

Our results show that unforeseen robustness is distinct from existing robustness metrics, further highlighting the need for a new measure which better captures the generalization of defense methods. We use ImageNet-UA to reveal that models with high L_∞ attack robustness (the most ubiquitous measure of robustness in the literature) do not generalize well to new attacks, recommending L_2 as a stronger baseline. We further find that L_p training can be improved on by alternative training processes, and suggest that the community focuses on training methods with better generalization behavior. Interestingly, unlike in the L_p case, we find that progress on CV benchmarks has at least partially tracked unforeseen robustness. We hope that the community can build on these results, so that ImageNet-UA can provide an improved progress measure for defenses aiming to achieve real-world worst-case robustness. To summarize, we make the following contributions:

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

a) Train against known adversaries

$$A_{\text{train}} = \left\{ \begin{array}{cc} \text{Image 1} & \text{Image 2} \\ L_2 & L_\infty \end{array} \right\} \text{ s.t. } A_{\text{train}} \cap A_{\text{test}} = \emptyset$$

b) Evaluate against our novel *unforeseen adversaries*

$$A_{\text{test}} = \left\{ \begin{array}{cccc} \text{Image 1} & \text{Image 2} & \text{Image 3} & \text{Image 4} \\ \text{Gabor} & \text{Elastic} & \text{Pixel} & \text{Glitch} \\ \text{Image 5} & \text{Image 6} & \text{Image 7} & \text{Image 8} \\ \text{Snow} & \text{JPEG} & \text{Wood} & \text{Kaleidoscope} \end{array} \right\}$$

Figure 1: **Testing unforeseen robustness.** Whereas most adversarial robustness benchmarks allow for training and evaluating on the same attacks, we evaluate robustness against a diverse range of held-out differentiable adversarial attacks. Thus, defences performing well on our ImageNet-UA demonstrate an ability to generalise to unforeseen adversaries.

- We design nineteen non- L_p attacks, constituting a large increase in the set of publicly-available non- L_p attacks. From this larger set, we carefully select eight core attacks for efficiency, efficacy and variety (see Appendix D for more discussion).
- We use our eight core attacks to form a new benchmark (ImageNet-UA), and a new metric UA2 (Unforeseen Adversarial Accuracy), standardizing and greatly expanding the scope of unforeseen robustness generalisation.
- We carry out in-depth analysis of the behaviour of different defence techniques on ImageNet-UA. We demonstrate that UA2 is distinct from existing robustness metrics in the literature, and point towards several promising directions for improving the generalisation of adversarial defences.

2 RELATED WORK

Evaluating Adversarial Robustness. Adversarial robustness is notoriously difficult to evaluate correctly (Papernot et al., 2017; Athalye et al., 2018). To this end, Carlini et al. (2019) provide extensive guidance for sound adversarial robustness evaluation. Our ImageNet-UA benchmark incorporates several of their recommendations, such as measuring attack success rates across several magnitudes of distortion and using a broader threat model with diverse differentiable attacks. Existing popular measures of adversarial robustness (Croce & Hein, 2020; Moosavi-Dezfooli et al., 2015; Weng et al., 2018) mainly present novel optimisation techniques for optimizing over an L_p -ball, limiting their applicability for modeling robustness to new deployment-time adversaries.

Non- L_p Attacks. Attacks often use hard-to-bound generative models (Song et al., 2018; Qiu et al., 2019), or make use of expensive brute-force search techniques Engstrom et al. (2017). We focus on attacks which are fast by virtue of differentiability, portable across datasets and independent of auxiliary generative models. Previous works presenting suitable attacks include Laidlaw & Feizi (2019); Shamsabadi et al. (2021); Zhao et al. (2019), who all transform the underlying color space of

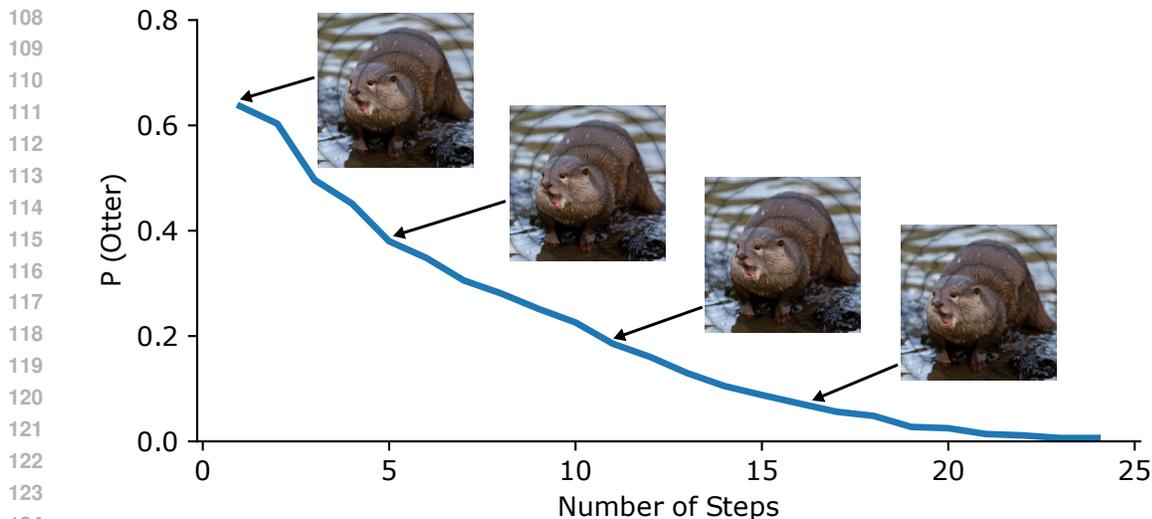


Figure 2: **Progression of an attack.** As we optimize our differentiable corruptions, model performance decreases, while leaving the image semantics unchanged. Unoptimized versions of our attacks have a moderate impact on classifier performance, similar to common corruptions (Hendrycks & Dietterich, 2019), while optimized versions are able to cause large drops in model accuracy.

an image and Xiao et al. (2018) who differentially warp images— a technique which we adapt to create our own Elastic adversary.

Unforeseen and Multi-attack Robustness. There exist defense methods which seek to generalize across an adversarial train-test gap (Dai et al., 2022; Laidlaw et al., 2020; Lin et al., 2020). Yet, comparison between these methods is challenging due to the lack of a standardized benchmark and an insufficient range of adversaries to test against. We fill this gap by implementing a unified benchmark for testing unforeseen robustness. The more developed field of multi-attack robustness (Tramer & Boneh, 2019) aims to create models which are robust to a range of attacks, but works generally focus on a union of L_p adversaries (Maini et al., 2020; Madaan et al., 2021a; Croce & Hein, 2022) and do not enforce that test time adversaries have to differ from those used during training.

Common corruptions Several of our attacks (Pixel, Snow, JPEG and Fog) were inspired by existing common corruptions (Hendrycks & Dietterich, 2018). We fundamentally change the generation methods to make these corruptions differentiable, allowing us to focus on worst-case robustness instead of the average-case robustness (see Section 4.1 for a comparison between the two benchmarks).

3 MEASURING UNFORESEEN ROBUSTNESS

To evaluate the unforeseen robustness of models, we introduce a new benchmark ImageNet-UA, and corresponding metric UA2 (Unforeseen Adversarial Accuracy). ImageNet-UA consists of our eight core adversarial attacks pictured in Figure 5, which have been selected for efficacy and computational efficiency.

3.1 THE UNFORESEEN ROBUSTNESS THREAT MODEL

The unforeseen robustness threat model has three central features:

Test-time distribution shift. As illustrated in Figure 5, our threat model requires that defenses do not make use of our test-time adversaries. This is directly analogous to the situation faced by developers, who are unable to anticipate which worst-case inputs which may occur at deployment time.

A diverse population of adversaries. Defence techniques should generalise to a range of possible deployment-time adversaries (e.g. they should not be effective only against L_p -based adversaries). Testing generalisation on a diverse range of adversaries is therefore an important part of evaluation.

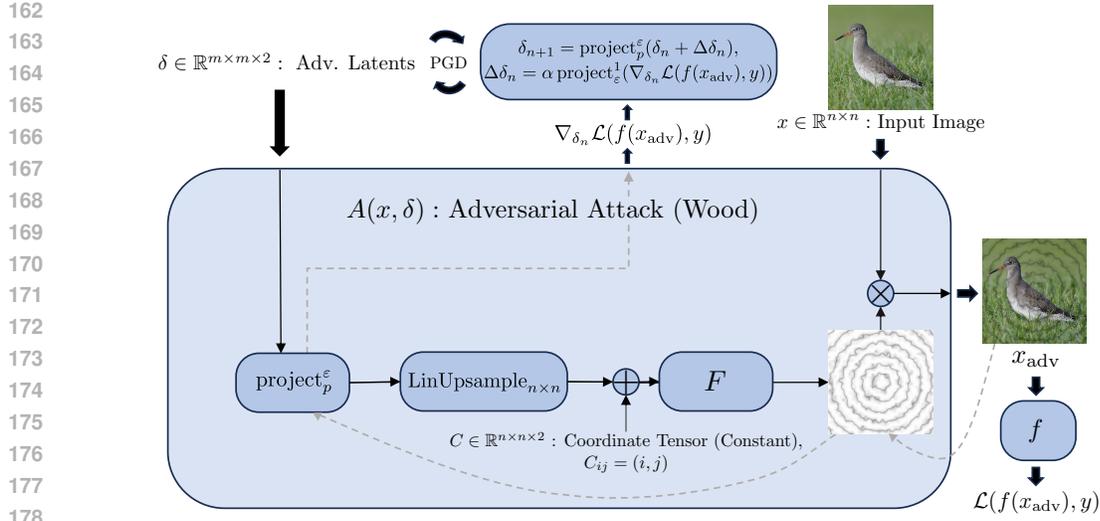


Figure 3: **An illustrative example of the generating process for one of our attacks.** As demonstrated by this illustration of our Wood attack, all of our attacks function by performing PGD optimization on a set of latent variables. In the case of the Wood attack, these latent variables are inputs to concentric sine waves ($F(x, y) = \sin(\sqrt{x^2 + y^2})$) which are overlaid on the image. See Appendix C for a more detailed explanation. We design effective attacks which are fast, easy to optimize, precisely bound, preserve image semantics, are portable across datasets and have variable intensity through the ε parameter.

White-box access to models. To ensure the strength of our adversaries (Carlini et al., 2019), and to avoid the usage of computationally expensive black-box optimization techniques, adversaries have full white-box access to victim models.

Unlike the more classical L_p threat models (Szegedy et al., 2013; Goodfellow et al., 2014a), we do not require that distortions are imperceptible (although they must preserve semantics, see Appendix K for a human study). This is because in many realistic settings, attackers are unlikely to face these constraints.

3.2 GENERATING ADVERSARIAL EXAMPLES

To test the performance of defence techniques, we must create a range of adversaries (see Section 3.3). Each of these adversaries are defined using a differentiable function A , which generates an adversarial input x_{adv} from an input image x and some latent variables δ :

$$x_{\text{adv}} = A(x, \delta). \quad (1)$$

To control the strength of our adversary, we introduce an L_p constraint to the variables δ , by bounding the size using some constraints ε_A :

$$\|\delta\|_p \leq \varepsilon_A$$

As is typical in the literature (Madry et al., 2017b), we use our dataset loss function \mathcal{L} to re-frame the finding of adversarial examples in our perturbation set as a continuous optimisation problem, seeking δ_{adv} which solves:

$$\delta_{\text{adv}} = \underset{\delta: \|\delta\|_p \leq \varepsilon}{\text{argmin}} \{ \mathcal{L}(f(A(x, \delta)), y) \}, \quad (2)$$

To solve this, we then use the popular method of Projected Gradient Descent (PGD) (Madry et al., 2017b) to find an approximate solution to Equation (2).

Using this formulation helps us ensure that all of our attacks are independent of auxiliary generative models, add minimal overhead when compared to the popular PGD adversary (see Appendix E), are usable in a dataset-agnostic “plug-and-play” manner, can be used with existing optimization algorithms (see Figure 4a for behavior of attacks under optimization), come with a natural way of

Table 1: L_p **robustness is distinct from unforeseen robustness**. We highlight some of the models which achieve high UA2, while still being susceptible to L_p attacks. Models below the dividing line are adversarially trained, with norm constraints in parentheses.

Model	L_∞ ($\varepsilon = 4/255$)	UA2
Dinov2 Vit-large	27.7	27.2
Convnext-V2-large IN-1k+22K	0.0	19.2
Swin-Large ImageNet1K	0.0	16.2
ConvNext-Base L_∞ , ($\varepsilon = 8/255$)	58.0	22.3
Resnet-50, L_∞ ($\varepsilon = 8/255$)	38.9	10
Resnet-50 L_2 , ($\varepsilon = 5$)	34.1	13.9

varying intensity through adjusting ε parameter (see Figure 4b for behavior under varying ε), and have precisely defined perturbation sets. As discussed in Section 2, this is not the case for most existing attacks in the literature, prompting us to design our set of new attacks.

3.3 EIGHT CORE UNFORESEEN ATTACKS

To ensure that we have high quality and diversity of tasks, we design nineteen novel attacks (see Appendix B for a full list), and select a core eight for their computational efficiency, effectiveness, variety and preservation of semantics (see Appendix D for further discussion, and Appendix K for a human study on semantic preservation). We use these attacks to create ImageNet-UA.

Importantly, we do not claim that these attacks provide an exhaustive taxonomy over adversaries. They simply represent a specific set of diverse held-out adversaries— mirroring methodologies that have been popular in other studies of distribution shift (Hendrycks & Dietterich, 2018; Hendrycks et al., 2021; Wang et al., 2019). We briefly describe our core attacks below:

Wood. The wood attack is described in Figure 3 and Appendix C.

Glitch. Glitch simulates a common behavior in corrupted images of colored fuzziness. Glitch greys out the image, splitting it into horizontal bars, before independently shifting color channels within each of these bars.

JPEG. The JPEG compression algorithm functions by encoding small image patches using the discrete cosine transform, and then quantizing the results. The attack functions by optimizing L_∞ -constrained perturbations within the JPEG-encoded space of compressed images and then reverse-transforming to obtain the image in pixel space, using ideas from Shin & Song (2017) to make this differentiable.

Gabor. Gabor spatially occludes the image with visually diverse Gabor noise (Lagae et al., 2009), optimizing the underlying sparse tensor which the Gabor kernels are applied to.

Kaleidoscope. Kaleidoscope overlays randomly colored polygons onto the image, and then optimizes both the homogeneous color of the inside of the shape, and the darkness/lightness of the individual pixels on the shape’s border, up to an L_∞ constraint.

Pixel. Pixel modifies an image so it appears to be of lower quality, by first splitting the image into $m \times m$ “pixels” and then averaging the image color within each block. The optimization variables δ then control the level of pixelation, on a per-block bases.

Elastic. Our only non-novel attack. Elastic is adapted from (Xiao et al., 2018), functioning by which warping the image by distortions $x' = \text{Flow}(x, V)$, where $V : \{1, \dots, 224\}^2 \rightarrow \mathcal{R}^2$ is a vector field on pixel space, and Flow sets the value of pixel (i, j) to the bilinearly interpolated original value at $(i, j) + V(i, j)$. To make the attack suitable for high-resolution images, we modify the original attack by passing a Gaussian kernel over V .

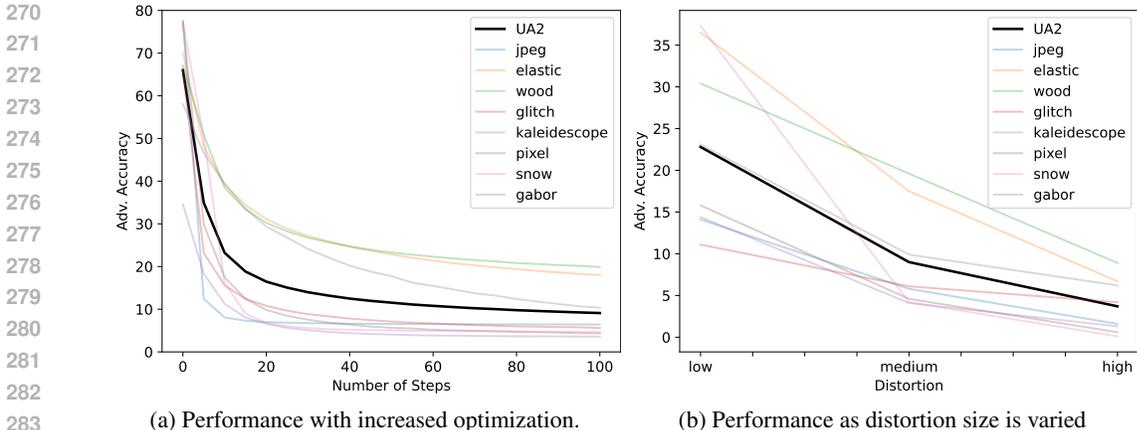


Figure 4: **Attack effectiveness increases with optimization pressure and distortion budget.** We average performance against our core attacks across all our benchmarked models, demonstrating that our attacks respond to increased optimization pressure (Figure 4a). We further demonstrate the importance of the gradient-based nature by comparing random grid search to our gradient-based method in Appendix M. Furthermore, we demonstrate the ability for our attack strength to be customisable by showing that increasing distortion budget reduces model performance (Figure 4b).

Snow. Snow functions by optimising the intensity of individually snowflakes within an image, which are created by passing a convolutional filter over a sparsely populated tensor, and then optimising the non-zero entries in this tensor.

Due to the time-consuming nature of designing new non- L_p adversarial attacks, most previous works only present and analyse a single adversarial attack (Laidlaw & Feizi, 2019; Shamsabadi et al., 2021; Zhao et al., 2019; Xiao et al., 2018). Hence, we believe that the additional attacks are a valuable resource for the community to build on, and release the full set of 19 publicly. We are particularly excited for these attacks as a basis of future evaluations of potential over-fitting to our original set of attacks, mirroring what was done by Mintun et al. (2021)

3.4 ImageNet-UA: A NEW BENCHMARK FOR UNFORESEEN ROBUSTNESS

We introduce ImageNet-UA, a benchmark for evaluating the unforeseen robustness of image classifiers on the popular ImageNet dataset (Deng et al., 2009). We also develop a CIFAR-10 equivalent, which we call CIFAR-10-UA, for computationally efficient evaluation of defense strategies and attack methods. We performed extensive sweeps to find the most effective hyperparameters for all of our attacks, the results of which are in Appendix A.

We quantify the unforeseen robustness achieved by a defense with Unforeseen Adversarial Accuracy (UA2), which measures the robustness of a given classifier f across a diverse range of attacks that are not seen at training time. We model the deployment-time population of adversaries as a uniform distribution over some finite set of adversaries \mathbf{A} . Our accuracy against such a population of adversaries is the average accuracy across each of the individual attacks:

$$UA2 := \frac{1}{|\mathbf{A}|} \sum_{A \in \mathbf{A}} Acc(A, \epsilon_A, f)$$

where $Acc(A, \epsilon_A, f)$ denotes the adversarial accuracy of classifier f against attack A at distortion level ϵ_A . We select the population of adversaries to be the eight core adversaries from Section 3.3, setting $\mathbf{A} = \{\text{JPEG}, \text{Elastic}, \text{Wood}, \text{Glitch}, \text{Kaleidoscope}, \text{Pixel}, \text{Snow}, \text{Gabor}\}$.

We further divide our benchmark by picking three different distortion levels for each attack, leading to three different measures of unforeseen robustness: $UA2_{\text{low}}$, $UA2_{\text{med}}$ and $UA2_{\text{high}}$ (see Appendix A for specific ϵ values used within this work), and we focus on focus on $UA2_{\text{med}}$ for all of our reports, referring to this distortion level as simply UA2. As distortion levels increase, model performance

Table 2: ImageNet-UA **baselines** Here, we show some of the most robust models on ImageNet-UA, as well as baseline ResNet-50 models to compare between. We see a variety of techniques achieving high levels of robustness, demonstrating a rich space of possible interventions. The L_∞ column tracks robustness against a PGD L_∞ adversary with $\varepsilon = 4/255$. Numbers denote percentages.

Model	Clean Acc.	L_∞	UA2	JPEG	Elastic	Wood	Glitch	Kal.	Pixel	Snow	Gabor
DINOv2 ViT-large Patch14	86.1	15.3	27.7	14.3	42.6	39.7	17.7	46.2	17.2	14.2	29.9
ConvNeXt-V2-large IN-1K+22K	87.3	0.0	19.2	0.0	39.1	34.4	21.4	16.1	15.5	4.0	23.1
ConvNeXt-V2-huge IN-1K	86.3	0.0	17.7	0.0	42.5	21.2	23.8	24.3	6.6	0.7	22.2
ConvNeXt-base, L_∞ (4/255)	76.1	58.0	22.3	39.0	23.8	47.9	12.9	2.5	9.7	30.2	12.8
ViT-base Patch16, L_∞ (4/255)	76.8	57.1	25.8	52.6	26.3	47.2	13.8	8.1	11.9	27.1	19.5
Swin-base IN-1K	85.3	0.0	15.2	0.0	31.4	24.6	16.2	6.0	6.9	4.3	32.0
ResNet-50	76.1	0.0	1.6	0.0	4.4	6.3	0.4	0.0	0.3	0.1	0.9
ResNet-50 + CutMix	78.6	0.5	6.1	0.2	17.9	15.5	2.5	0.1	6.7	3.0	2.7
ResNet-50, L_∞ (8/255)	54.5	38.9	10.0	6.9	11.8	23.9	14.4	0.7	5.2	15.6	1.2
ResNet-50, L_2 (5)	56.1	34.1	13.9	39.7	11.9	19.4	12.2	0.3	9.7	15.4	2.5

Table 3: L_p **training**. We train a range of ResNet-50 models against L_p adversaries on ImageNet-UA

Training	Train ε	Clean Acc.	UA2
Standard	-	76.1	1.6
L_2	1	69.1	6.4
	3	62.8	12.2
	5	56.1	13.9
L_∞	2/255	69.1	6.4
	4/255	63.9	7.9
	8/255	54.5	10.0

Table 4: L_p **training on generated data**. We see the effect of training when training WRN-28-10 networks on CIFAR-10-50M, a 1000x larger diffusion-model generated version of CIFAR-10 (Wang et al., 2023)

Dataset	Training	Clean Acc.	UA2
CIFAR-10	$L_2, \varepsilon = 1$	82.3	45.8
	$L_\infty, \varepsilon = 8/255$	86.1	41.5
CIFAR-10-50M	$L_2, \varepsilon = 0.5$	95.2	51.2
	$L_\infty, \varepsilon = 4/255$	92.4	51.5

decreases (Figure 4b). We perform a human study (Appendix K) to ensure UA2_{med} preserves image semantics.

4 EXPERIMENTS

In this section, we evaluate a range of models on our standardized benchmarks ImageNet-UA and CIFAR-10-UA. We aim to present a set of directions for future work, by comparing a wide range of methods. We also hope to explore how the problem of unforeseen robustness differs from existing robustness metrics.

4.1 HOW DO EXISTING ROBUSTNESS MEASURES RELATE TO UNFORESEEN ROBUSTNESS?

We find differences between existing metrics and UA2, suggesting that the setting of unforeseen adversarial robustness may require new methods to obtain strong performance.

UA2 is distinct from existing measures of distribution shift. We compare UA2 to several standard distribution-shift benchmarks—ImageNet-C (Hendrycks & Dietterich, 2019), ImageNet-R Hendrycks et al. (2021) and ImageNet-Sketch (Wang et al., 2019). As shown in Table 5 and Appendix I, performance on these benchmark is similar to performance on non-optimized versions of our attacks. By contrast, the optimized versions of our attacks are far more challenging and have distinct properties. For example, while adversarial training does not improve performance on ImageNet-C, it does improve performance on ImageNet-UA. This highlights that UA2 is a measure of worst case robustness, similar to L_p robustness, and distinct from other distribution shift benchmarks in the literature.

L_p robustness is correlated, but distinct from, unforeseen robustness. As shown in Appendix L, unforeseen robustness is correlated with L_p robustness. Our attacks also show similar properties to L_p counterparts, such as the ability for black-box transfer (Appendix N). However, many models show susceptibility to L_p adversaries while still performing well on UA2 (Table 1), and a range

Table 5: **Common corruptions and UA2** We compare ImageNet-C to both non-optimized and optimized versions of our attacks. We find that ImageNet-C behaves similarly to our non-optimized attacks, while our optimised attacks are far more challenging.

Model	UA2 (non-optimized) \uparrow	mCE \downarrow	UA2 \uparrow
Resnet 50	55.2	76.7	1.6
Resnet50 + AugMix	59.1	65.7	3.5
Resnet50 + DeepAug	60.2	61.1	3.0
Resnet50 + Mixup	59.9	69.2	4.8
Resnet50 + L_2 , ($\varepsilon = 5$)	43.2	89.0	13.9
Resnet50 + L_∞ , ($\varepsilon = 8/255$)	40.6	85.1	10

of strategies beat L_p training baselines Section 4.2 . We conclude that UA2 is distinct from L_p robustness, and present UA2 as an improved progress measure when working towards real-world worst-case robustness.

L_2 -based adversarial training outperforms L_∞ -based adversarial training. We see that L_p adversarial training increases the unforeseen robustness of tested models, with L_2 adversarial training providing the largest increase in UA2 over standard training (1.6% \rightarrow 13.9%), beating models which are trained against L_∞ adversaries (1.6% \rightarrow 10.0%). We present L_2 trained models as a strong baseline for unforeseen robustness, noting that the discrepancy between L_∞ and L_2 training is particularly relevant as L_∞ robustness is the most ubiquitous measure of adversarial robustness.

4.2 HOW CAN WE IMPROVE UNFORESEEN ROBUSTNESS?

We find several promising directions that improve over standard L_p training. These include standard multi-attack robustness methods and also novel methods that combine insights from different areas of robustness research.

Combining image augmentations and L_∞ training. One simple approach to improving robustness to unforeseen adversaries would be to combine insights from distribution shift robustness and adversarial robustness research. To explore this avenue, we experiment with a combination of PixMix and L_∞ adversarial training, applying adversarial perturbations to the augmented images from PixMix. We show the results of this experiment in Table 6.

Surprisingly, we find that this simple approach can be highly effective. Namely, we find that combining PixMix and adversarial training gives a UA2 of 45.5 percent, compared to 37.3 with adversarial training alone. This novel training strategy beats strong baselines by combining two distinct robustness techniques. The surprising effectiveness of this simple method highlights how unforeseen robustness may foster the development of new methods.

Multi-attack robustness. To evaluate how existing work on robustness to a union of L_p balls may improve unforeseen robustness, we use CIFAR-10-UA to evaluate a strong multi-attack robustness baseline by (Madaan et al., 2021b), which trains a Meta Noise Generator (MNG) that learns the optimal training perturbations to achieve robustness to a union of L_p adversaries. For WRN-28-10 models on CIFAR-10-UA, we see a large increase in unforeseen robustness compared to the best L_p baseline (21.4% \rightarrow 51.1%, full results in Appendix J), leaving scaling of such methods to full ImageNet-UA for future work.

Bounding perturbations with perceptual distance. We evaluate the UA2 of models trained with Perceptual Adversarial Training (PAT) (Laidlaw et al., 2020). PAT functions by training a

Table 6: **PixMix and L_p training.** We compare UA2 on CIFAR-10 of models trained with PixMix and adversarial training. Combining PixMix and adversarial training improves UA2, demonstrating the potential for novel methods to improve UA2. All numbers denote percentages, and L_∞ training was performed with the TRADES algorithm.

Training Strategy	Train ε	Clean Acc.	UA2
PixMix	-	95.1	15.00
L_∞	4/255	89.3	37.3
L_∞ + PixMix	4/255	91.4	45.1
L_∞	8/255	84.3	41.4
L_∞ + PixMix	8/255	87.1	47.4

Table 7: **Effects of data augmentation on UA2.** We evaluate the UA2 of a range of data-augmented ResNet50 models.

Training	Clean Acc.	UA2
Standard	76.1	1.0
Moex	79.1	6.0
CutMix	78.6	6.0
Deepaugment + Augmix	75.8	1.8

Table 8: **Effects of pretraining and regularization on UA2.**

Model	Clean Acc.	UA2
ConvNeXt-V2-28.6M	83.0	9.8
ConvNeXt-V1-28M	82.1	5.1
ConvNeXt-V2-89M	84.9	14.9
ConvNeXt-V1-89M	83.8	9.7
ConvNeXt-V2-198M	85.8	19.1
ConvNeXt-V1-198M	84.3	10.6

model against an adversary bounded by an estimate of the human perceptual distance, computing the estimate by using the hidden states of an image classifier. For computational reasons we train and evaluate ResNet-50s on a 100-image subset of ImageNet-UA, where this technique outperforms the best L_p trained baselines (22.6 \rightarrow 26.2, full results in Appendix J).

Regularizing high-level features. We evaluate Variational Regularization (VR) (Dai et al., 2022), which adds a penalty term to the loss function for variance in higher level features. We find that the largest gains in unforeseen robustness come from combining VR with PAT, improving over standard PAT (26.2 \rightarrow 29.5, on a 100 class subset of ImageNet-UA, full results in Appendix J).

4.3 HOW HAS PROGRESS ON CV BENCHMARKS TRACKED UNFORESEEN ROBUSTNESS?

Computer vision progress has partially tracked unforeseen robustness. Comparing the UA2 of ResNet-50 to ConvNeXt-V2-huge (1% \rightarrow 19.1% UA2) demonstrates the effects of almost a decade of CV advances, including self-supervised pretraining, hardware improvements, data augmentation, architectural changes and new regularization techniques. More generally, we find a range of modern architectures and training strategies doing well (see Table 2, full results in Figure 8). This is gives a positive view of how progress on standard CV benchmarks has tracked underlying robustness metrics, contrasting with classical L_p adversarial robustness where standard training techniques have little effect (Madry et al., 2017a).

Scale, data augmentation and pretraining successfully improve unforeseen robustness. We do a more careful analysis of how three of the most effective CV techniques have improved robustness. As shown in Section 4.2, we find that data augmentation improves on unforeseen robustness, even in cases where they reduce standard accuracy. We compare the performance of ConvNeXt-V1 and ConvNeXt-V2 models, which differ through the introduction of self-supervised pretraining and a new normalization layer. When controlling for model capacity, the combination of this pre-training and normalisation layer demonstrates a large increase unforeseen robustness Table 8.

5 CONCLUSION

In this paper, we introduced a new benchmark for testing robustness against *unforeseen adversaries* (ImageNet-UA) laying groundwork for continuing research in improving real world adversarial robustness. We provide nineteen (eighteen novel) non- L_p attacks as part of our repository, using these to construct a new metric UA2 (Unforeseen Adversarial Accuracy). We show that ImageNet-UA is distinct from existing measures of robustness in the literature, and make use use it to evaluate classical L_p training techniques—showing that the common practice of L_∞ training and evaluation may be misleading, as L_2 training shows higher unforeseen robustness. We additionally demonstrate that a variety of interventions outside of L_p adversarial training can improve unforeseen robustness, both through existing techniques in the CV literature and through more specialised adversarial training strategies. We hope that ImageNet-UA will be a useful tool as we continue to make progress towards safer machine learning systems in real-world applications.

REFERENCES

- 486
487
488 Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of
489 security: Circumventing defenses to adversarial examples. In *International conference on machine*
490 *learning*, pp. 274–283. PMLR, 2018.
- 491 Tom B. Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch.
492 *CoRR*, abs/1712.09665, 2017. URL <http://arxiv.org/abs/1712.09665>.
- 493 John Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis*
494 *and machine intelligence*, (6):679–698, 1986.
- 496 Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras,
497 Ian J. Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness.
498 *CoRR*, abs/1902.06705, 2019. URL <http://arxiv.org/abs/1902.06705>.
- 499 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble
500 of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–2216.
501 PMLR, 2020.
- 502
503 Francesco Croce and Matthias Hein. Adversarial robustness against multiple and single l_p -threat
504 models via quick fine-tuning of robust classifiers, 2022.
- 505 Sihui Dai, Saeed Mahloujifar, and Prateek Mittal. Formulating robustness against unforeseen attacks.
506 *arXiv preprint arXiv:2204.13779*, 2022.
- 507
508 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
509 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
510 pp. 248–255. Ieee, 2009.
- 511 Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. A
512 rotation and a translation suffice: Fooling CNNs with simple transformations. *arXiv preprint*
513 *arXiv:1712.02779*, 2017.
- 514
515 Alain Fournier, Don Fussell, and Loren Carpenter. Computer rendering of stochastic models.
516 *Commun. ACM*, 25(6):371–384, June 1982. ISSN 0001-0782. doi: 10.1145/358523.358553. URL
517 <http://doi.acm.org/10.1145/358523.358553>.
- 518 Justin Gilmer, Ryan P. Adams, Ian J. Goodfellow, David Andersen, and George E. Dahl. Motivating
519 the rules of the game for adversarial example research. *ArXiv*, abs/1807.06732, 2018.
- 520
521 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
522 examples. *arXiv preprint arXiv:1412.6572*, 2014a.
- 523
524 Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
525 examples, 2014b. URL <https://arxiv.org/abs/1412.6572>.
- 526 Qing Guo, Felix Juefei-Xu, Xiaofei Xie, Lei Ma, Jian Wang, Bing Yu, Wei Feng, and
527 Yang Liu. Watch out! motion is blurring the vision of your deep neural networks. In
528 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in*
529 *Neural Information Processing Systems*, volume 33, pp. 975–985. Curran Associates, Inc.,
530 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/0a73de68f10e15626eb98701ecf03adb-Paper.pdf)
531 [file/0a73de68f10e15626eb98701ecf03adb-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/0a73de68f10e15626eb98701ecf03adb-Paper.pdf).
- 532 Qing Guo, Ziyi Cheng, Felix Juefei-Xu, Lei Ma, Xiaofei Xie, Yang Liu, and Jianjun Zhao. Learning to
533 adversarially blur visual object tracking. In *Proceedings of the IEEE/CVF International Conference*
534 *on Computer Vision*, pp. 10839–10848, 2021.
- 535 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
536 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*
537 *Computer Vision and Pattern Recognition*, pp. 16000–16009, 2022.
- 538
539 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
corruptions and perturbations. In *International Conference on Learning Representations*, 2019.

- 540 Dan Hendrycks and Thomas G Dietterich. Benchmarking neural network robustness to common
541 corruptions and surface variations. *arXiv preprint arXiv:1807.01697*, 2018.
- 542
- 543 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
544 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical
545 analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International
546 Conference on Computer Vision*, pp. 8340–8349, 2021.
- 547 Ares Lagae, Sylvain Lefebvre, George Drettakis, and Philip Dutré. Procedural noise using sparse
548 Gabor convolution. *ACM Trans. Graph.*, 28(3):54:1–54:10, July 2009. ISSN 0730-0301. doi:
549 10.1145/1531326.1531360. URL <http://doi.acm.org/10.1145/1531326.1531360>.
- 550 Cassidy Laidlaw and Soheil Feizi. Functional adversarial attacks. *Advances in neural information
551 processing systems*, 32, 2019.
- 552
- 553 Cassidy Laidlaw, Sahil Singla, and Soheil Feizi. Perceptual adversarial robustness: Defense against
554 unseen threat models, 2020. URL <https://arxiv.org/abs/2006.12655>.
- 555 Wei-An Lin, Chun Pong Lau, Alexander Levine, Rama Chellappa, and Soheil Feizi. Dual manifold
556 adversarial robustness: Defense against lp and non-lp adversarial attacks. *Advances in Neural
557 Information Processing Systems*, 33:3487–3498, 2020.
- 558
- 559 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
560 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the
561 IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 562 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.
563 A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
564 Pattern Recognition*, pp. 11976–11986, 2022.
- 565 Divyam Madaan, Jinwoo Shin, and Sung Ju Hwang. Learning to generate noise for multi-attack
566 robustness, 2021a.
- 567
- 568 Divyam Madaan, Jinwoo Shin, and Sung Ju Hwang. Learning to generate noise for multi-attack
569 robustness. In *International Conference on Machine Learning*, pp. 7279–7289. PMLR, 2021b.
- 570 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
571 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,
572 2017a.
- 573
- 574 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
575 Towards deep learning models resistant to adversarial attacks, 2017b. URL [https://arxiv.
576 org/abs/1706.06083](https://arxiv.org/abs/1706.06083).
- 577 Pratyush Maini, Eric Wong, and J. Zico Kolter. Adversarial robustness against the union of multiple
578 perturbation models, 2020.
- 579 Karttikeya Mangalam, Haoqi Fan, Yanghao Li, Chao-Yuan Wu, Bo Xiong, Christoph Feichtenhofer,
580 and Jitendra Malik. Reversible vision transformers. In *Proceedings of the IEEE/CVF Conference
581 on Computer Vision and Pattern Recognition*, pp. 10830–10840, 2022.
- 582
- 583 Eric Mintun, Alexander Kirillov, and Saining Xie. On interaction between augmentations and
584 corruptions in natural corruption robustness. *Advances in Neural Information Processing Systems*,
585 34:3571–3583, 2021.
- 586 Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. DeepFool: a simple and
587 accurate method to fool deep neural networks. *arXiv preprint arXiv:1511.04599*, 2015.
- 588
- 589 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
590 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
591 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 592 Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram
593 Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on
Asia conference on computer and communications security*, pp. 506–519. ACM, 2017.

- 594 Ken Perlin. Making noise, 1999. URL <http://www.noisemachine.com/talk1/index.html>, 2005.
- 595
- 596 Vinay Uday Prabhu. The blood diamond effect in neural art : On ethically troublesome images of the
597 imagenet dataset. 2019.
- 598 Haonan Qiu, Chaowei Xiao, Lei Yang, Xinchun Yan, Honglak Lee, and Bo Li. Semanticadv:
599 Generating adversarial examples via attribute-conditional image editing. *ArXiv*, abs/1906.07927,
600 2019.
- 601
- 602 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
603 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
604 models from natural language supervision. In *International conference on machine learning*, pp.
605 8748–8763. PMLR, 2021.
- 606 Leslie Rice, Eric Wong, and Zico Kolter. Overfitting in adversarially robust deep learning. In
607 *International Conference on Machine Learning*, pp. 8093–8104. PMLR, 2020.
- 608
- 609 Ali Shahin Shamsabadi, Changjae Oh, and Andrea Cavallaro. Semantically adversarial learnable
610 filters. *IEEE Transactions on Image Processing*, 30:8075–8087, 2021.
- 611
- 612 Richard Shin and Dawn Song. JPEG-resistant adversarial images. In *NIPS 2017 Workshop on
613 Machine Learning and Computer Security*, 2017.
- 614 Yang Song, Rui Shu, Nate Kushman, and Stefano Ermon. Constructing unrestricted adversarial
615 examples with generative models. In *NeurIPS*, 2018.
- 616
- 617 Andreas Peter Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit,
618 and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision
619 transformers. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL
620 <https://openreview.net/forum?id=4nPswr1KcP>.
- 621 David Stutz, Matthias Hein, and Bernt Schiele. Confidence-calibrated adversarial training: Generaliz-
622 ing to unseen attacks. In *International Conference on Machine Learning*, pp. 9155–9166. PMLR,
623 2020.
- 624
- 625 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,
626 and Rob Fergus. Intriguing properties of neural networks. 12 2013.
- 627 Florian Tramèr and Dan Boneh. Adversarial training and robustness for multiple perturbations.
628 *Advances in neural information processing systems*, 32, 2019.
- 629
- 630 Florian Tramèr, Pascal Dupré, Gili Rusak, Giancarlo Pellegrino, and Dan Boneh. Ad-versarial:
631 Defeating perceptual ad-blocking. *CoRR*, abs/1811.03194, 2018. URL [http://arxiv.org/
632 abs/1811.03194](http://arxiv.org/abs/1811.03194).
- 633 Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations
634 by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32,
635 2019.
- 636
- 637 Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion
638 models further improve adversarial training, 2023.
- 639 Tsui-Wei Weng, Huan Zhang, Pin-Yu Chen, Jinfeng Yi, Dong Su, Yupeng Gao, Cho-Jui Hsieh, and
640 Luca Daniel. Evaluating the robustness of neural networks: An extreme value theory approach.
641 *arXiv preprint arXiv:1801.10578*, 2018.
- 642
- 643 Sanghyun Woo, Shoubhik Debnath, Ronghang Hu, Xinlei Chen, Zhuang Liu, In So Kweon, and
644 Saining Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. *arXiv
645 preprint arXiv:2301.00808*, 2023.
- 646
- 647 Chaowei Xiao, Jun-Yan Zhu, Bo Li, Warren He, Mingyan Liu, and Dawn Song. Spatially trans-
formed adversarial examples. In *Proceedings of (ICLR) International Conference on Learning
Representations*, April 2018.

648 Yaodong Yu, Zitong Yang, Edgar Dobriban, Jacob Steinhardt, and Yi Ma. Understanding generaliza-
649 tion in adversarial training via the bias-variance decomposition. *arXiv preprint arXiv:2103.09947*,
650 2021.

651
652 Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
653 Theoretically principled trade-off between robustness and accuracy. In Kamalika Chaudhuri
654 and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine*
655 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 7472–7482, Long Beach,
656 California, USA, 09–15 Jun 2019. PMLR. URL [http://proceedings.mlr.press/v97/
657 zhang19p.html](http://proceedings.mlr.press/v97/zhang19p.html).

658 Zhengyu Zhao, Zhuoran Liu, and Marisa Larson. Towards large yet imperceptible adversarial image
659 perturbations with perceptual color distance. *ArXiv*, abs/1911.02466, 2019.

660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A HYPERPARAMETERS

A.1 TRAINED MODELS

To run our evaluations, we train a range of our own models to benchmark with:

- CIFAR-10 WRN-28-10 robust models and TRADES models are respectively trained with the official code of Rice et al. (2020) and Zhang et al. (2019) with the default hyperparameters settings
- The PAT-VR models on ImageNet100 were trained using the official code from Dai et al. (2022) and employed the hyperparameter settings outlined in the code of Laidlaw et al. (2020).
- ImageNet100 DINOv2 Oquab et al. (2023) models are trained by finetuning a linear classification head on the ImageNet100 dataset. We used a SGD optimizer with learning rate of 0.001 and employed early-stopping.

A.2 MODEL REFERENCE

We use a range of baseline models provided by other works, with model weights available as part of their open source distribution:

- **ImageNet**
 - ConvNeXt models are from Liu et al. (2022)
 - ConvNeXt-V2 models are from Woo et al. (2023)
 - ViT models are from Steiner et al. (2022)
 - Swin models are from Liu et al. (2021)
 - Reversible-ViT models are from Mangalam et al. (2022)
 - CLIP (ViT-L/14) is from Radford et al. (2021)
 - DINOv2 models are from Oquab et al. (2023)
 - MAE models are from He et al. (2022)
- **CIFAR-10**
 - WideResNet TRADES models are from Zhang et al. (2019)
 - WRN + Diffusion models are from Wang et al. (2023)
 - Meta noise models are from Madaan et al. (2021b)
 - ResNet50 VR models are from Dai et al. (2022)
 - ReColorAdv models are from Laidlaw & Feizi (2019)
 - StAdv models are from Xiao et al. (2018)
 - Multi attack models are from Tramèr et al. (2018)
 - The Multi steepest descent model is from Maini et al. (2020)
 - PAT models are from Laidlaw et al. (2020)
 - Pre-trained ResNet18 L_∞ , L_2 and L_1 models are from Croce & Hein (2022)
- **ImageNet100**
 - ResNet50 PAT models are from Laidlaw et al. (2020)
 - ResNet50 PAT + VR models are from Dai et al. (2022)
 - DINOv2 models are from Oquab et al. (2023)

A.3 ATTACK PARAMETERS

To ensure that our attacks are maximally effective, we perform extensive hyper-parameter sweeps to find the most effective step sizes.

Table 9: Attack parameters for ImageNet-UA

		Step Size	Num Steps	Low Distortion	Medium Distortion	High Distortion	Distance Metric
	PGD	0.004	50	2/255	4/255	8/255	L_∞
	Gabor	0.0025	100	0.02	0.04	0.06	L_∞
	Snow	0.1	100	10	15	25	L_2
	Pixel	1	100	3	5	10	L_2
Core Attacks	JPEG	0.0024	80	1/255	3/255	6/255	L_∞
	Elastic	0.003	100	0.1	0.25	0.5	L_2
	Wood	0.005	80	0.03	0.05	0.1	L_∞
	Glitch	0.005	90	0.03	0.05	0.07	L_∞
	Kaleidoscope	0.005	90	0.05	0.1	0.15	L_∞
	Edge	0.02	60	0.03	0.1	0.3	L_∞
	FBM	0.006	30	0.03	0.06	0.3	L_∞
	Fog	0.05	80	0.3	0.5	0.7	L_∞
	HSV	0.012	50	0.01	0.03	0.05	L_∞
	Klotski	0.01	50	0.03	0.1	0.2	L_∞
Extra Attacks	Mix	1.0	70	5	10	40	L_2
	Pokadot	0.3	70	1	3	5	L_2
	Prison	0.0015	30	0.01	0.03	0.1	L_∞
	Blur	0.03	40	0.1	0.3	0.6	L_∞
	Texture	0.00075	80	0.01	0.03	0.2	L_∞
	Whirlpool	4.0	40	10	40	100	L_2

Table 10: Attack parameters for CIFAR-10-UA

		Step Size	Num Steps	Low Distortion	Medium Distortion	High Distortion	Distance Metric
	PGD	0.008	50	2/255	4/255	8/255	L_∞
	Gabor	0.0025	80	0.02	0.03	0.04	L_∞
	Snow	0.2	20	3	4	5	L_2
	Pixel	1.0	60	1	5	10	L_2
Core Attacks	JPEG	0.0024	50	1/255	3/255	6/255	L_∞
	Elastic	0.006	30	0.1	0.25	0.5	L_2
	Wood	0.000625	70	0.03	0.05	0.1	L_∞
	Glitch	0.0025	60	0.03	0.05	0.1	L_∞
	Kaleidoscope	0.005	30	0.05	0.1	0.15	L_∞
	Edge	0.02	60	0.03	0.1	0.3	L_∞
	FBM	0.006	30	0.02	0.04	0.08	L_∞
	Fog	0.05	40	0.3	0.4	0.5	L_∞
	HSV	0.003	20	0.01	0.02	0.03	L_∞
	Klotski	0.005	50	0.03	0.05	0.1	L_∞
Extra Attacks	Mix	0.5	30	1	5	10	L_2
	Pokadot	0.3	40	1	2	3	L_2
	Prison	0.0015	20	0.01	0.03	0.1	L_∞
	Blur	0.015	20	0.1	0.3	0.6	L_∞
	Texture	0.003	30	0.01	0.1	0.2	L_∞
	Whirlpool	16.0	50	20	100	200	L_2

B DESCRIPTIONS OF THE 11 ADDITIONAL ATTACKS.

Blur. Blur approximates real-world motion blur effects by passing a Gaussian filter over the original image and then does a pixel-wise linear interpolation between the blurred version and the original, with the optimisation variables controlling the level of interpolation. We also apply a Gaussian filter to the grid of optimisation variables, to enforce some continuity in the strength of the blur between adjacent pixels. This method is distinct from, but related to other blurring attacks in the literature (Guo et al., 2020; 2021).

Edge. This attack functions by applying a Canny Edge Detector (Canny, 1986) over the image to locate pixels at the edge of objects, and then applies a standard PGD attack to the identified edge pixels.

Fractional Brownian Motion (FBM). FBM overlays several layers of Perlin noise (Perlin, 2005) at different frequencies, creating a distinctive noise pattern. The underlying gradient vectors which generate each instance of the Perlin noise are then optimised by the attack.

Fog. Fog simulates worst-case weather conditions, creating fog-like occlusions by adversarially optimizing parameters in the diamond-square algorithm (Fournier et al., 1982) typically used to render stochastic fog effects.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

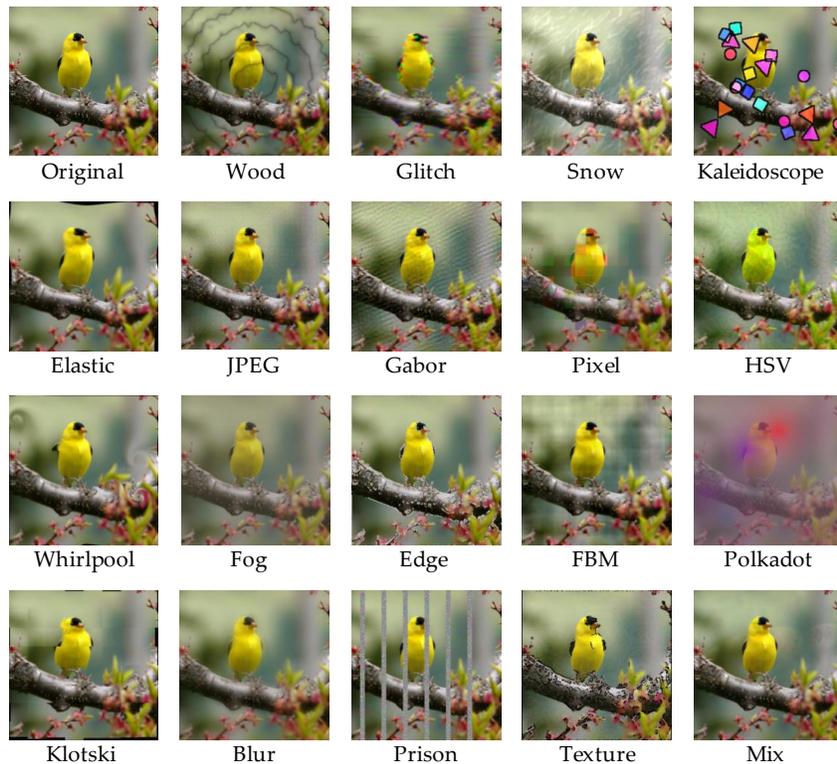


Figure 5: **The full suite of attacks.** We present nineteen differentiable non-Lp attacks as part of our codebase. For the purpose of visualization, higher distortion levels than are used in our benchmark have been chosen. See Appendix G for adversarial examples generated with the distortion levels used within our benchmark, and Appendix K for a human study on semantic preservation

HSV. This attack transforms the image into the HSV color space, and then optimises PGD in that latent space. Due to improving optimisation properties, a gaussian filter is passed over the image.

Klotski. The Klotski attack works by splitting the image into blocks, and applying a differentiable translation to each block, which is then optimised.

Mix. The Mix attack functions by performing differentiable pixel-wise interpolation between the original image and an image of a different class. The level of interpolation at each pixel is optimised, and a gaussian filter is passed over the pixel interpolation matrix to ensure that the interpolation is locally smooth.

Polkadot. Polkadot randomly selects points on the image to be the centers of a randomly coloured circle, and then optimising the size of these circles in a differentiable manner.

Prison. Prison places grey “prison bars” across the image, optimising only the images within the prison bars. This attack is inspired by previous “patch” attacks (Brown et al., 2017), while ensuring that only the prison bars are optimised.

Texture. Texture works by removing texture information within an images, passing a Canny Edge Detector (Canny, 1986) over the image to find all the pixels which are at the edges of objects, and then filling these pixels in black—creating a silhouette of the original image. The other non-edge (or “texture”) pixels are then whitened, losing the textural information of the image while preserving the shape. Per-pixel optimisation variable control the level of whitening.

Whirlpool. Whirlpool translates individual pixels in the image by a differentiable function creating a whirlpool-like warping of the image, optimising the strength of each individual whirlpool.

C FULL DESCRIPTION OF WOOD ATTACK

In Figure 3, we give a high-level explanation of the Wood attack. Here, we give a more detailed explanation of this figure.

Given a classifier f , the Wood attack with distortion level ε functions by taking a set of adversarial latent variables $\delta_n \in \mathbb{R}^{m \times m \times 2}$ (representing a vector field of per-pixel displacements), applies $project_p^\varepsilon$ to project this field into the ε ball in the L_p metric and then uses bi-linear interpolation to upsample the latent variables to the input image size. The upsampled latent variables are then used to make the wood noise, by using an element-wise mapping $F: \mathbb{R}^{n \times n \times 2} \rightarrow \mathbb{R}^{n \times n}$, taking a coordinate to the (power of) the sine of its distance from the center of the image i.e. $F(I) = \sin(\sqrt{(X)^2 + (Y)^2})^\beta$, where $X_{ij} = I_{ij0} - n/2$ and $Y_{ij} = I_{ij1} - n/2$ and β is an attack hyperparameter. When applied to constant coordinate tensor $C \in \mathbb{R}^{n \times n \times 2}$, $C_{ij} = (i, j)$, this function creates the distinctive “wood rings” of the Wood attack, which are then multiplied with the input image to produce adversarial input. By virtue of the differentiability of this process, we can backpropagate through this noise generation and optimize the adversarial image x_{adv} by performing PGD (Madry et al., 2017a) on the input latent variables.

D PROCESS FOR DESIGNING ATTACKS AND SELECTING CORE ATTACKS

Our design of attacks is guided by two motivations: defending against unforeseen adversaries and robustness to long-tail scenarios. Unforeseen adversaries could implement novel attacks to, e.g., evade automated neural network content filters. To model unforeseen adversaries that might realistically appear in these scenarios, we include digital corruptions similar to what one might see on YouTube videos trying to evade content filters. These include attacks such as Kaleidoscope and Prison. To model long-tail scenarios, we include worst-case versions of common corruptions, like JPEG, Snow, and Fog.

In preliminary experiments, we found that some of these attacks were more effective than others, leading to lower accuracy with fewer steps. We also found that performance on some attacks was correlated between models. For example, both “Prison” and “Edge”, are pixel-level attacks, so robustness to one was correlated with robustness to the other. To increase the diversity and efficiency of our evaluation, we selected a core set of eight attacks based on their effectiveness and diversity, considering both visual diversity and the accuracy profiles of different models. This was an iterative process that led us to make substantial changes to some attacks. For example, we modified the implementation of the Elastic attack to use larger, lower-frequency distortions, which maintained its effectiveness while reducing correlation with PGD.

E ATTACK COMPUTATION TIME

We investigate the execution times of our attacks, finding that most attacks are not significantly slower than an equivalent PGD adversary.

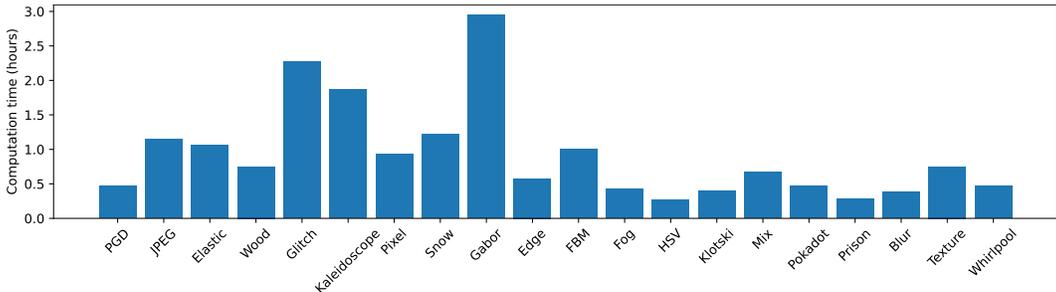


Figure 6: Evaluation time of the attacks on the ImageNet test set using a ResNet50 model with batch size of 200 on a single A100-80GB GPU, Attack hyper-parameters are as described in Appendix A.

F FULL RESULTS OF MODEL EVALUATIONS

We benchmark a large variety of models on our dataset, finding a rich space of interventions affecting unforeseen robustness.

F.1 IMAGENET

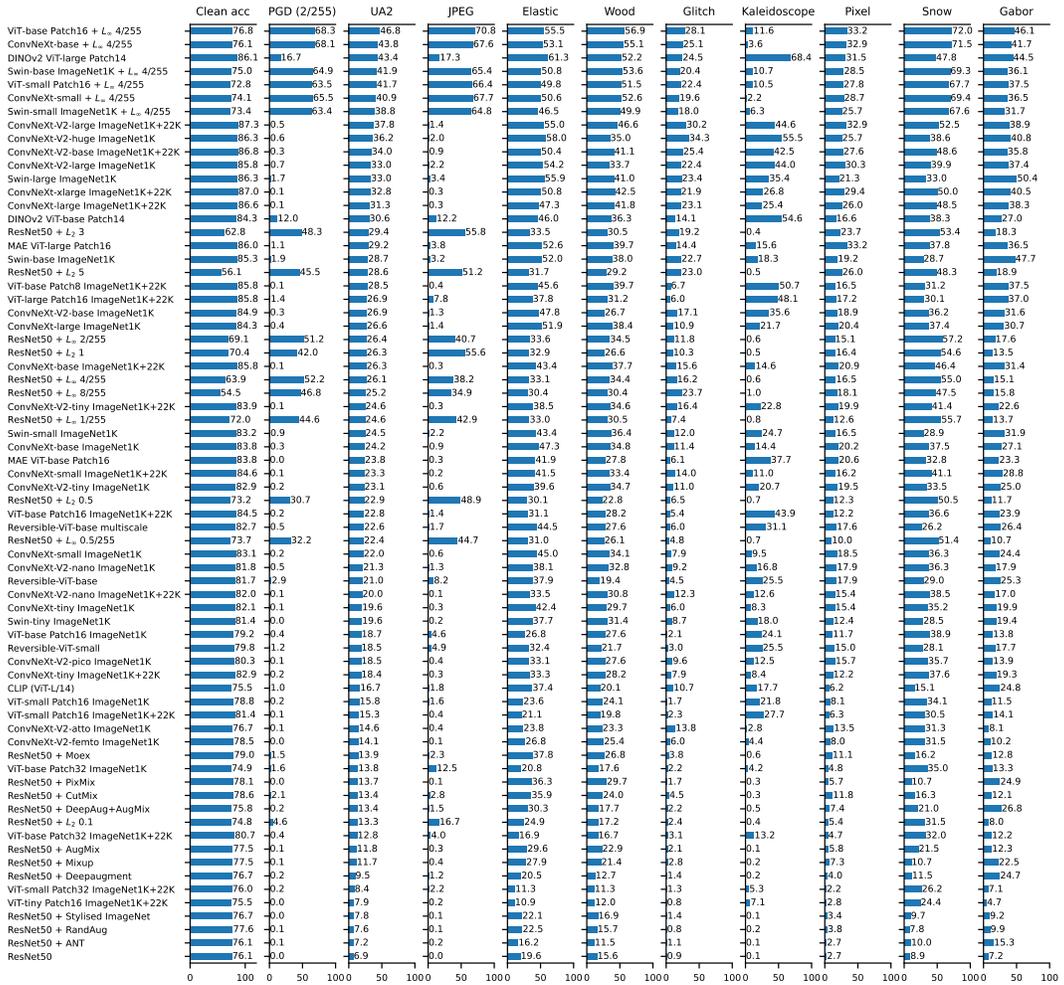


Figure 7: ImageNet UA2 performance under low distortion.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

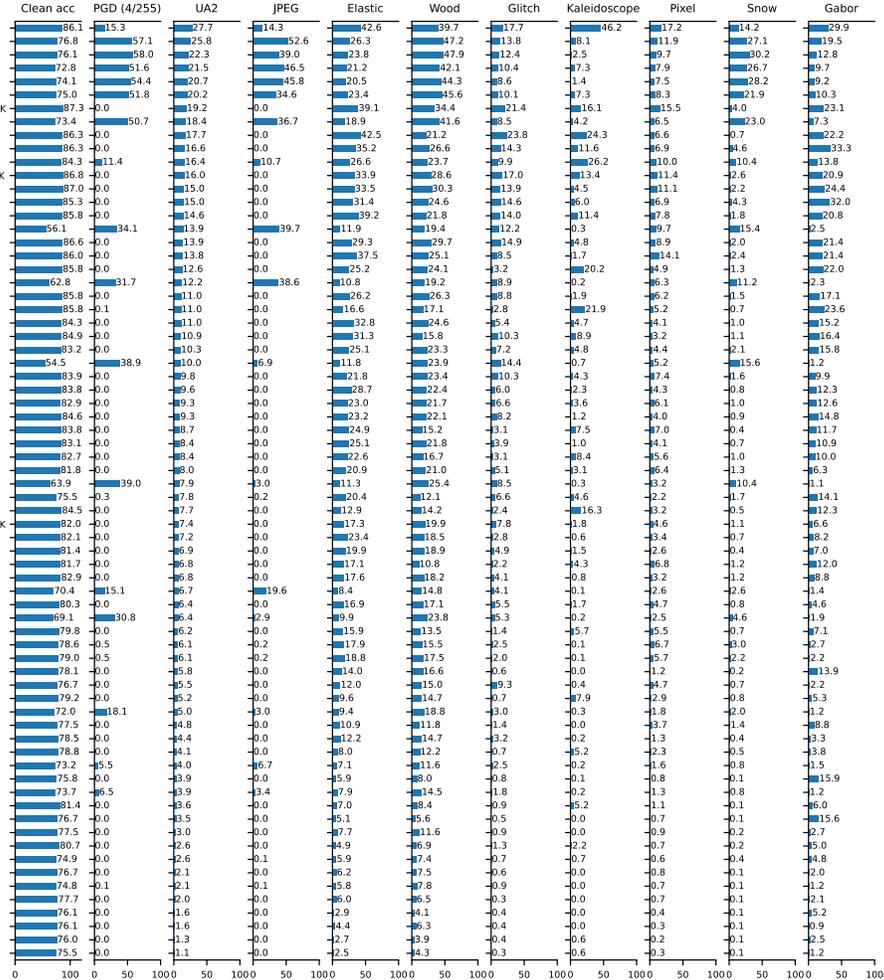


Figure 8: ImageNet UA2 performance under medium distortion

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

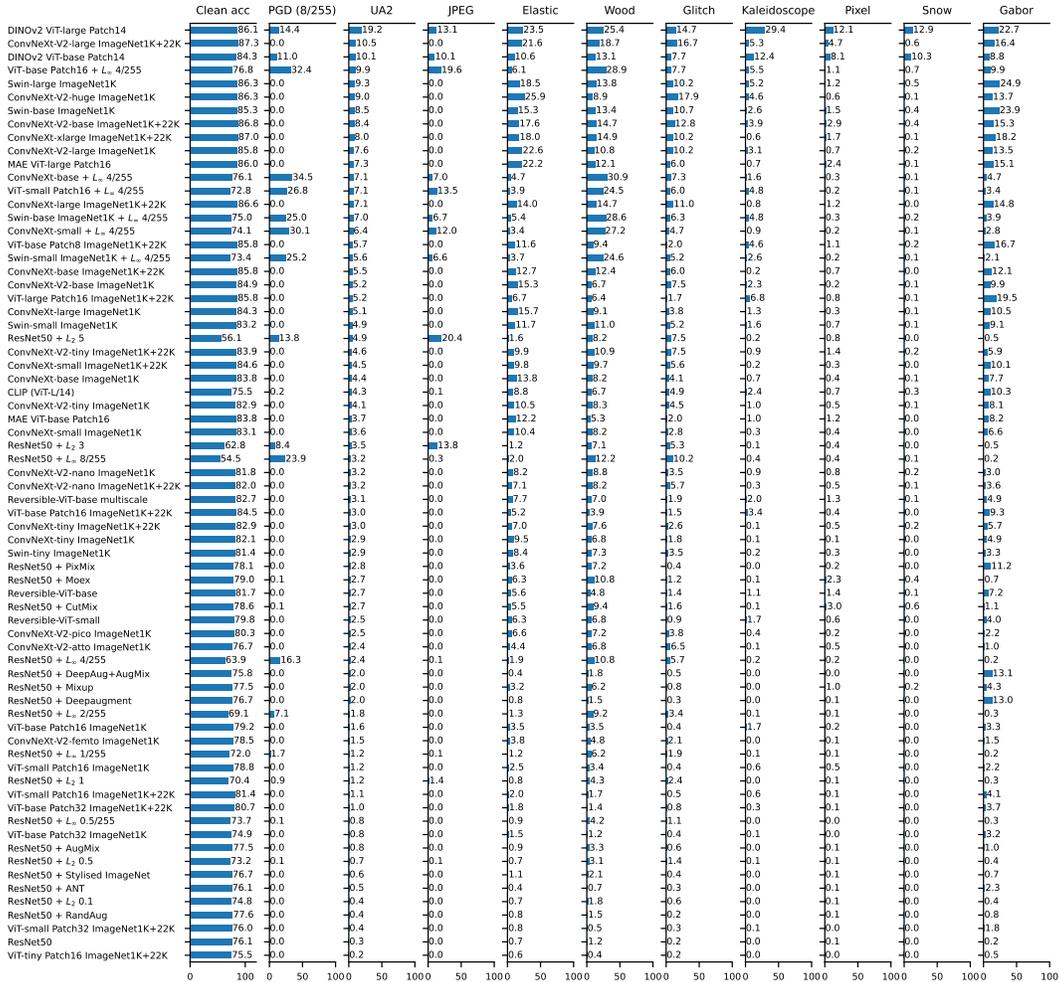


Figure 9: ImageNet UA2 performance under high distortion

F.2 CIFAR-10

	Clean acc	PGD (2/255)	UA2	JPEG	Elastic	Wood	Glitch	Kaleidoscope	Pixel	Snow	Gabor
WRN-70-16 + L ₂ 128/255 + Diffusion model	95.5	71.3	75.6	92.6	85.8	63.0	83.0	42.4	90.8	83.2	63.8
WRN-28-10 + L ₂ 128/255 + Diffusion model	95.2	90.2	74.1	91.8	84.1	61.3	81.8	40.4	89.7	81.4	62.6
WRN-70-16 + L ₂ 8/255 + Diffusion model	93.3	90.1	73.3	87.3	83.5	55.4	78.8	46.6	84.6	80.0	69.8
WRN-28-10 + L ₂ 8/255 + Diffusion model	92.4	88.9	72.4	86.1	81.9	53.9	77.5	46.6	83.4	80.1	69.3
WRN-28-10 + Meta noise + Robust self-training	88.9	83.0	70.1	84.7	76.0	61.5	79.0	43.3	81.7	73.0	61.9
WRN-28-10 + L ₂ 8/255 + Pixmix	91.4	84.5	68.8	82.5	80.6	55.9	75.7	33.4	79.1	71.3	61.5
WRN-28-10 + L ₂ 8/255 + Pixmix	87.1	81.5	67.1	79.4	76.8	55.9	74.1	36.2	74.9	71.7	68.1
ResNet50 + Multi attack (average)	85.9	78.9	66.0	79.5	77.8	58.2	75.0	36.2	76.2	66.8	64.0
ResNet50 + Multi attack (maximum)	83.5	77.4	64.7	71.9	77.7	61.2	70.3	36.8	71.8	67.7	60.5
WRN-28-10 + Meta noise	81.5	74.5	64.2	77.1	68.0	57.5	71.7	37.5	73.4	61.9	63.1
WRN-28-10 + L ₂ 0.5 (TRADES)	87.6	78.9	64.2	82.2	71.2	53.5	75.0	33.1	78.0	60.8	59.4
WRN-28-10 + Multi attack (maximum)	84.3	76.7	64.2	78.2	70.0	57.2	72.3	37.3	74.8	61.3	62.2
PreAct ResNet18 Union of L _p	80.5	74.1	64.0	76.2	67.2	56.3	70.4	42.4	72.4	63.4	63.4
WRN-28-10 + Multi attack (average)	85.8	77.0	63.7	79.9	67.9	56.4	73.0	33.8	76.3	60.4	61.6
WRN-28-10 + L ₂ 1.0 (TRADES)	82.3	75.2	63.5	78.3	67.7	58.0	72.2	35.1	74.0	61.0	61.7
WRN-34-10 + L ₂ 8/255	86.1	80.6	63.3	78.1	70.0	44.1	73.6	38.8	74.0	65.7	61.9
WRN-28-10 + Multi steepest descent	82.7	76.1	63.2	77.0	68.0	56.7	71.5	36.5	72.8	61.5	61.7
WRN-34-10 + L ₂ 8/255 (TRADES)	84.9	79.4	63.2	76.4	71.8	44.0	72.1	42.6	71.9	65.7	61.2
ResNet50 + L ₂ 0.5	89.6	79.9	63.1	83.4	70.9	48.7	69.4	28.6	78.9	63.0	62.0
PreAct ResNet18 L ₁	80.7	72.8	63.1	76.1	65.8	54.3	71.1	39.5	73.3	62.3	62.2
ResNet18 + L ₂ 0.5 + VR 1.0	85.2	76.2	62.8	79.6	64.6	49.7	72.5	35.7	75.2	60.1	64.6
ResNet50 + PAT 0.5 + VR 0.1	94.8	87.0	62.7	89.4	71.0	43.9	67.4	32.6	86.6	64.6	61.0
WRN-28-10 + L ₂ 0.5 (TRADES)	91.7	80.4	62.6	83.8	74.9	43.8	74.8	30.2	80.3	59.5	53.6
ResNet50 + PAT 0.5 + VR 0.05	96.3	77.6	62.6	80.5	71.7	44.6	73.3	29.9	76.7	62.4	61.4
WRN-28-10 + L ₂ 4/255	91.9	84.6	62.5	82.1	75.5	38.6	74.4	32.0	78.0	64.3	55.2
WRN-28-10 + L ₂ 4/255 (TRADES)	89.3	82.2	62.4	79.1	71.6	38.1	74.0	36.6	74.8	63.0	58.7
ResNet50 + PAT 0.5	95.7	76.5	62.3	79.8	70.5	48.3	72.4	28.6	75.4	61.1	62.1
WRN-28-10 + L ₂ 8/255	86.5	80.5	62.2	77.9	70.3	42.4	74.2	34.9	73.7	66.0	58.5
WRN-28-10 + L ₂ 8/255 (TRADES)	84.3	78.0	61.6	75.0	70.1	42.0	72.1	37.8	71.3	63.8	60.5
ResNet50 + L ₂ 8/255	95.5	79.5	61.3	77.4	71.6	45.6	63.4	30.5	72.2	66.4	63.7
ResNet50 + L ₂ 0.25	92.0	79.7	61.0	83.0	75.0	35.4	68.5	30.0	80.3	59.7	56.0
ResNet50 + Self-bounded PAT	82.1	74.1	60.8	76.7	65.6	47.2	71.7	35.9	74.3	59.4	55.4
ResNet50 + L ₂ 1.0	79.0	72.2	59.8	74.8	66.2	61.9	60.8	31.6	70.9	62.5	60.0
WRN-40-2 + L ₂ 8/255	93.3	86.5	63.5	73.9	66.1	39.7	62.8	35.0	85.7	60.3	61.5
ResNet50 + Multi attack (random)	81.8	70.2	57.8	73.5	68.1	50.3	63.3	28.8	66.8	55.7	56.1
ResNet50 + ReColorAdv	93.0	74.2	57.6	79.4	74.8	30.1	68.0	20.1	79.4	63.0	46.3
PreAct ResNet18 L ₁ pretrained	87.2	68.0	56.4	75.1	64.0	35.2	67.2	36.3	72.7	48.4	52.0
ResNet50 + AlexNet-bounded PAT	71.1	67.0	55.9	69.6	62.0	48.6	63.2	28.2	66.2	64.9	55.4
ResNet18 + L ₂ 8/255 + VR 0.5	72.9	68.4	55.8	66.2	55.5	38.6	63.9	41.9	62.1	60.2	58.2
ResNet50 + PAT 1.0 + VR 0.05	71.4	67.0	55.8	68.4	62.1	43.6	63.5	33.7	66.0	66.0	52.8
WRN-28-10 + L ₂ 8/255 + VR 0.7	72.7	68.8	55.7	66.0	56.9	39.0	63.7	42.4	61.9	59.3	56.3
ResNet50 + PAT 1.0 + VR 0.1	71.5	66.2	55.7	69.0	59.7	44.8	62.2	36.0	65.3	64.9	54.4
WRN-28-10 + L ₂ 0.5	95.6	74.9	55.5	79.6	74.0	23.0	60.4	29.3	80.6	55.2	42.1
WRN-28-10 + L ₂ 1.0	95.5	75.3	55.4	79.6	73.7	22.2	61.0	30.3	80.2	55.1	40.8
PreAct ResNet18 L ₁ pretrained	82.9	70.9	54.9	66.7	63.6	26.5	66.5	38.9	66.8	53.6	56.4
WRN-28-10 + L ₂ 0.25	95.3	74.0	54.6	78.6	73.3	21.8	59.8	27.4	79.8	53.9	42.1
ResNet18 + ReColorAdv + VR 1.0	94.0	66.1	53.8	75.3	70.9	20.6	59.3	23.9	77.9	54.6	47.6
ResNet18 + ReColorAdv + VR 0.5	94.0	65.5	53.1	75.1	71.7	19.3	59.8	21.2	78.0	54.0	46.1
ResNet50 + SAdv	85.7	64.1	50.3	64.1	60.2	65.3	39.6	37.8	45.5	71.2	8.9
ResNet18 + SAdv + VR 1.0	90.6	66.4	50.3	66.9	74.9	52.8	44.7	41.5	68.5	67.1	15.9
ResNet18 + SAdv	86.8	63.3	49.8	62.4	61.2	61.9	39.3	33.6	45.5	71.2	13.0
ResNet18 + SAdv + VR 0.5	82.9	68.9	49.5	61.1	76.2	59.9	39.3	39.6	44.6	67.2	18.0
WRN-40-2 + Augmix	95.0	76.0	46.2	33.9	67.5	17.9	45.5	31.9	71.4	61.7	49.4
ResNet18 + ReColorAdv	94.7	69.5	46.1	63.5	65.8	61.1	49.7	19.9	76.9	44.4	40.9
PreAct ResNet18 L ₁ pretrained	81.5	46.7	44.6	49.0	50.6	27.3	59.3	25.9	65.3	42.1	37.6
WRN-40-2 + Pixmix	95.1	20.6	39.6	14.7	54.0	7.9	37.9	19.1	70.1	47.9	65.1
WRN-40-2 + Uniform noise	94.3	21.8	30.4	25.0	46.5	2.5	22.2	23.3	64.9	30.3	28.2
ResNet50	94.7	6.5	29.2	12.5	47.2	1.7	17.0	14.1	64.6	26.9	39.6
WRN-28-10	95.8	10.4	28.5	13.3	43.5	2.6	19.2	17.9	64.7	38.3	28.7
WRN-40-2 + Cutmix	95.7	5.4	26.6	5.2	38.2	2.9	13.4	30.4	60.2	40.3	22.4
WRN-40-2 + Standard	94.6	6.4	24.3	6.3	33.1	1.1	13.9	16.2	62.7	28.4	30.4
WRN-40-2 + Mixup	94.8	6.2	24.2	5.0	32.7	2.1	14.5	13.9	60.7	36.0	28.5

Figure 10: CIFAR-10 UA2 performance under low distortion.

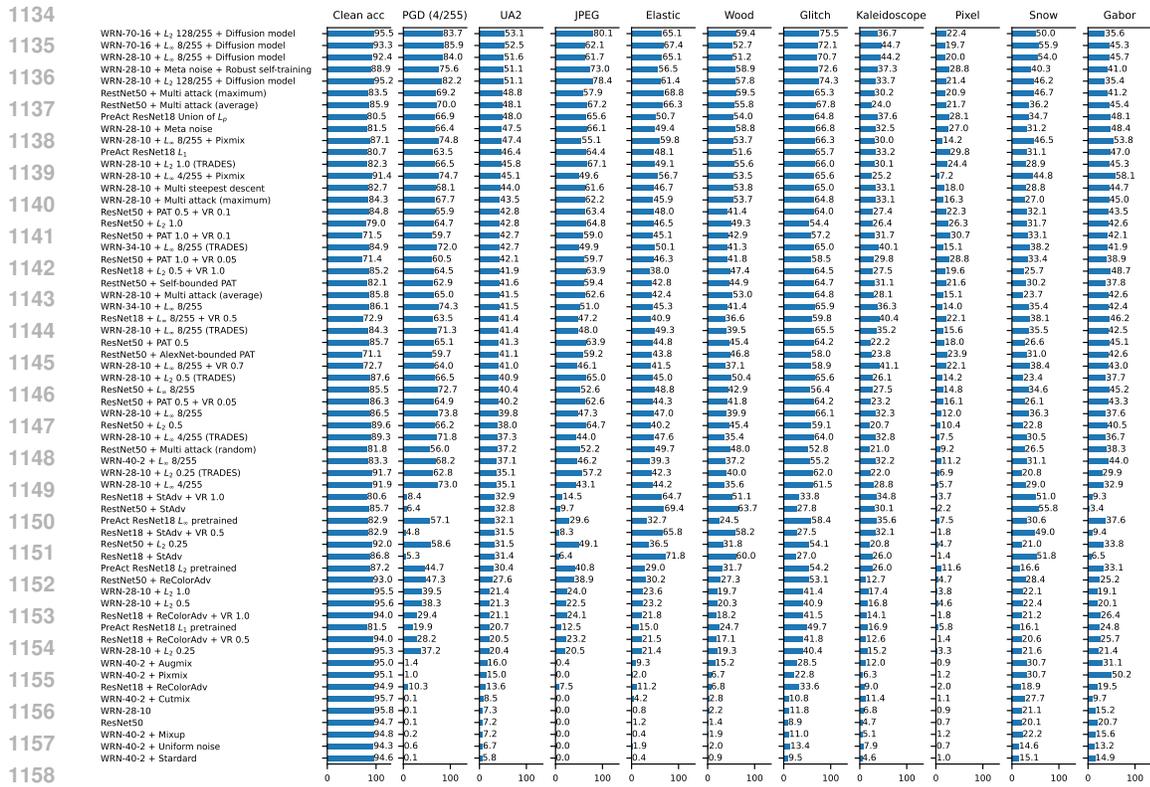
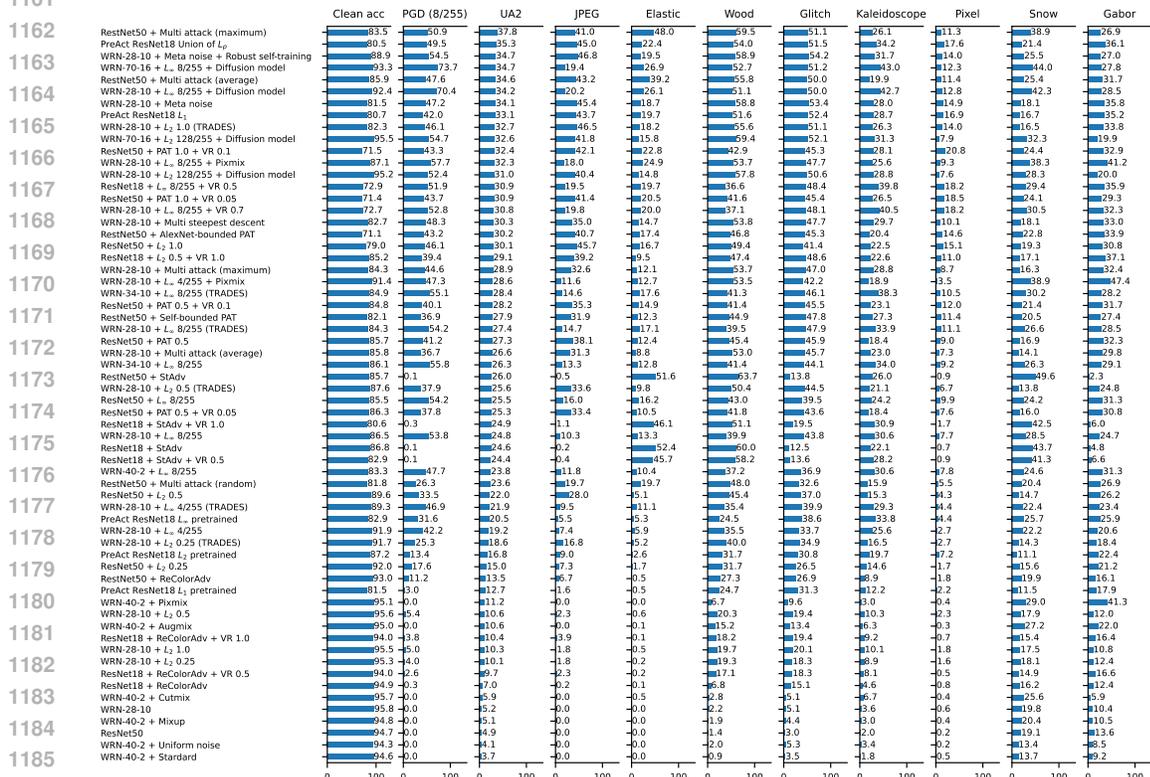


Figure 11: CIFAR-10 UA2 performance under medium distortion



F.3 IMAGENET100

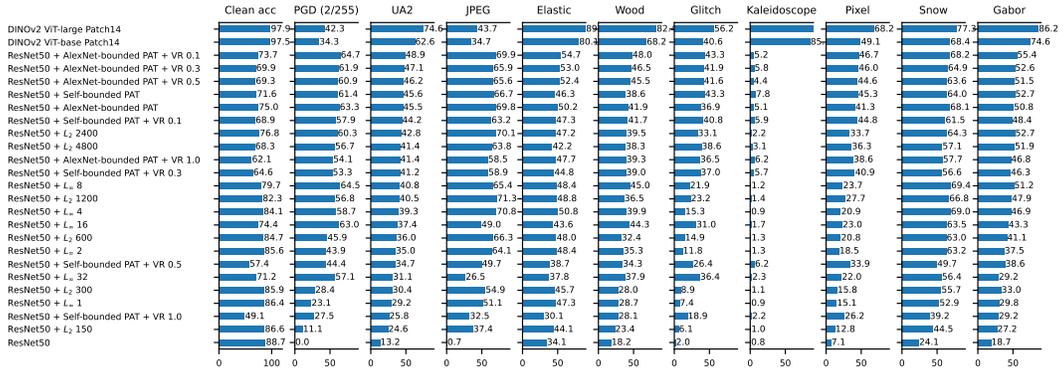


Figure 13: ImageNet100 UA2 performance under low distortion.

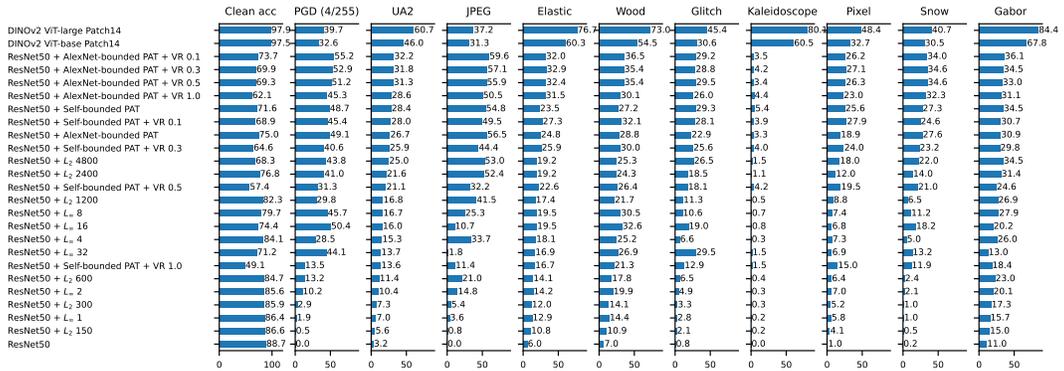


Figure 14: ImageNet100 UA2 performance under medium distortion

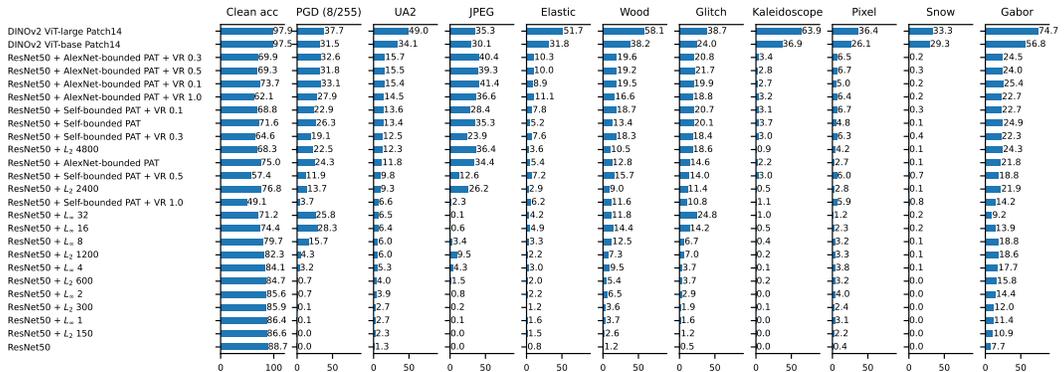


Figure 15: ImageNet100 UA2 performance under high distortion

F.4 EXPLORING THE ROBUSTNESS OF DINOv2

Given the strong adversarial robustness of DINOv2 models under the PGD attack (Appendix F), we further evaluate the DINOv2 model under AutoAttack Croce & Hein (2020). Table 11 and Table 12 show that although for the robust ResNet50 model AutoAttack performs similarly to PGD, it is able to reduce the accuracy of DINOv2 models to 0.0% across all the distortion levels. Future work may benefit from applying the AutoAttack benchmark as a comparison point, instead of the base PGD adversary.

Table 11: Attacked accuracies of models on ImageNet

	ResNet50 + L_∞ 8/255	DINOv2 ViT-base Patch14	DINOv2 ViT-large Patch14
PGD (2/255)	46.8%	12.0%	16.7%
APGD-CE (2/255)	46.2%	1.0%	1.0%
APGD-CE + APGD-T (2/255)	43.6%	0.0%	0.0%
PGD (4/255)	38.9%	11.4%	15.3%
APGD-CE (4/255)	37.9%	0.9%	0.8%
APGD-CE + APGD-T (4/255)	33.8%	0.0%	0.0%
PGD (8/255)	23.9%	11.0%	14.4%
APGD-CE (8/255)	22.6%	0.6%	0.7%
APGD-CE + APGD-T (8/255)	18.4%	0.0%	0.0%

Table 12: Attacked accuracies of models on ImageNet100

	ResNet50 + L_∞ 8/255	DINOv2 ViT-base Patch14	DINOv2 ViT-large Patch14
PGD (2/255)	64.5%	34.3%	42.3%
APGD-CE (2/255)	64.4%	17.6%	20.0%
APGD-CE + APGD-T (2/255)	64.1%	0.0%	0.0%
PGD (4/255)	45.7%	32.6%	39.7%
APGD-CE (4/255)	45.2%	16.4%	17.3%
APGD-CE + APGD-T (4/255)	44.6%	0.0%	0.0%
PGD (8/255)	15.7%	31.5%	37.7%
APGD-CE (8/255)	14.7%	15.5%	14.5%
APGD-CE + APGD-T (8/255)	13.6%	0.0%	0.0%

F.5 PERFORMANCE VARIANCE

As described in Section 3.2, we perform adversarial attacks by optimizing latent variables which are randomly initialized in our current implementation, so the adversarial attack’s performance can be affected by the random seed for the initialization. To study the effect of random initializations, we compute the UA2 performances of three samples of two ImageNet models, ResNet50 and ResNet50 + L_2 5. We observe the standard deviations of UA2 of these two models across 5 different seeds to be respectively 0.1% and 0.04% concluding that the variation in performance across the ImageNet dataset is minor.

G IMAGES OF ALL ATTACKS ACROSS DISTORTION LEVELS

We provide images of all 19 attacks within the benchmark, across the three distortion levels.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

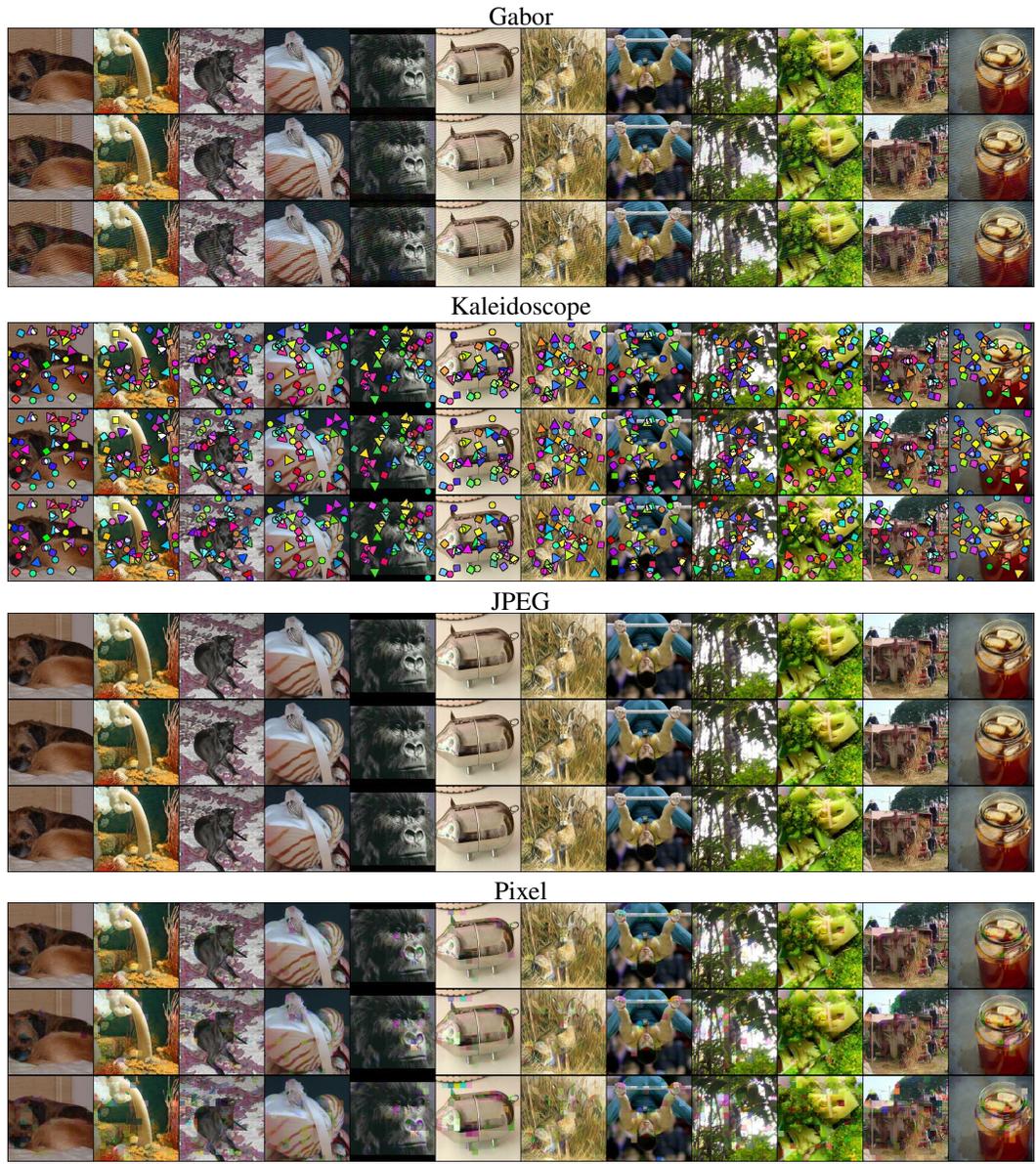


Figure 16: Attacked samples of low distortion (1st row), medium distortion (2nd row), and high distortion (last row) on a standard ResNet50 model

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403



Figure 17: Attacked samples of low distortion (1st row), medium distortion (2nd row), and high distortion (last row) on a standard ResNet50 model

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

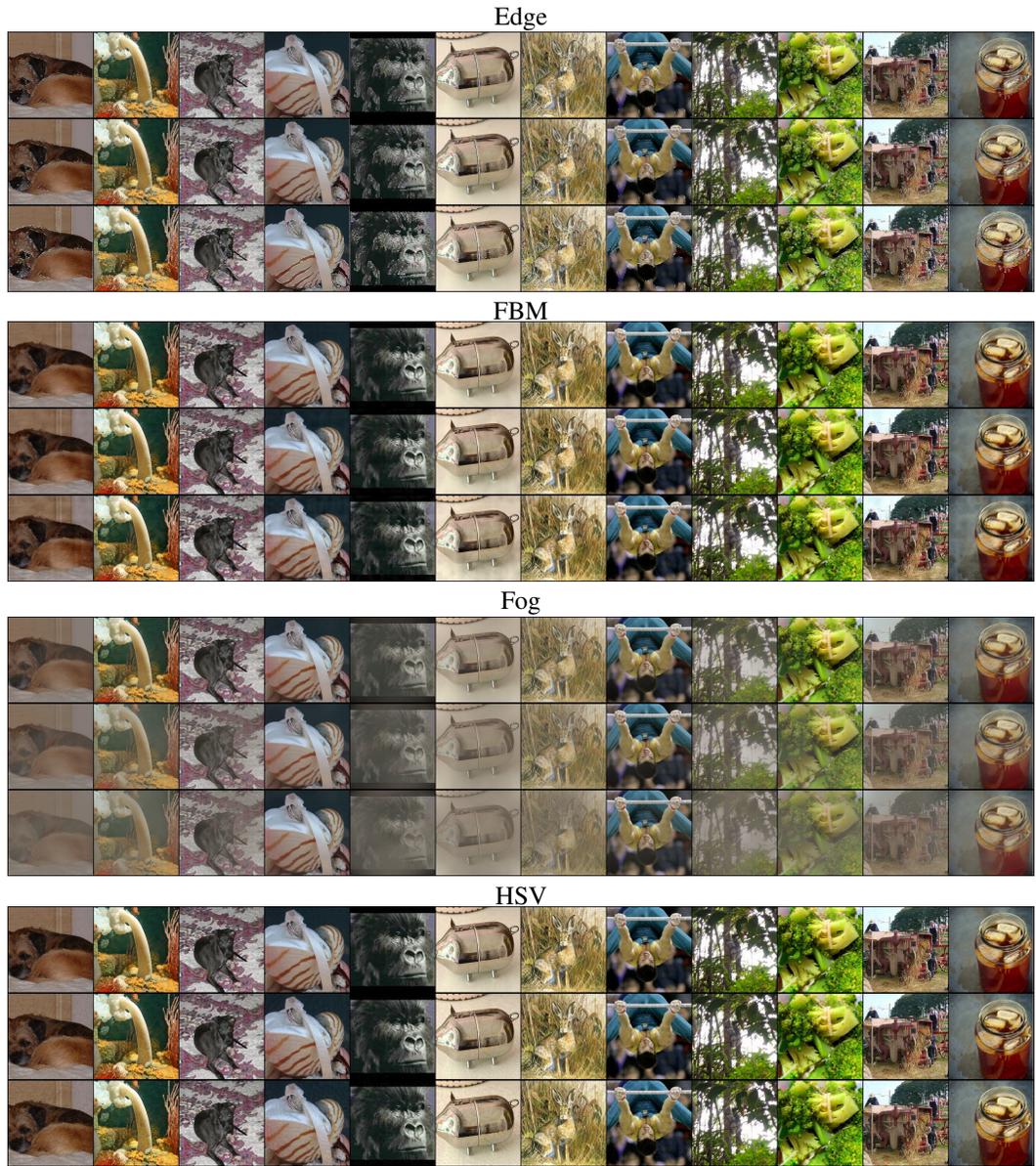


Figure 18: Attacked samples of low distortion (1st row), medium distortion (2nd row), and high distortion (last row) on a standard ResNet50 model

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

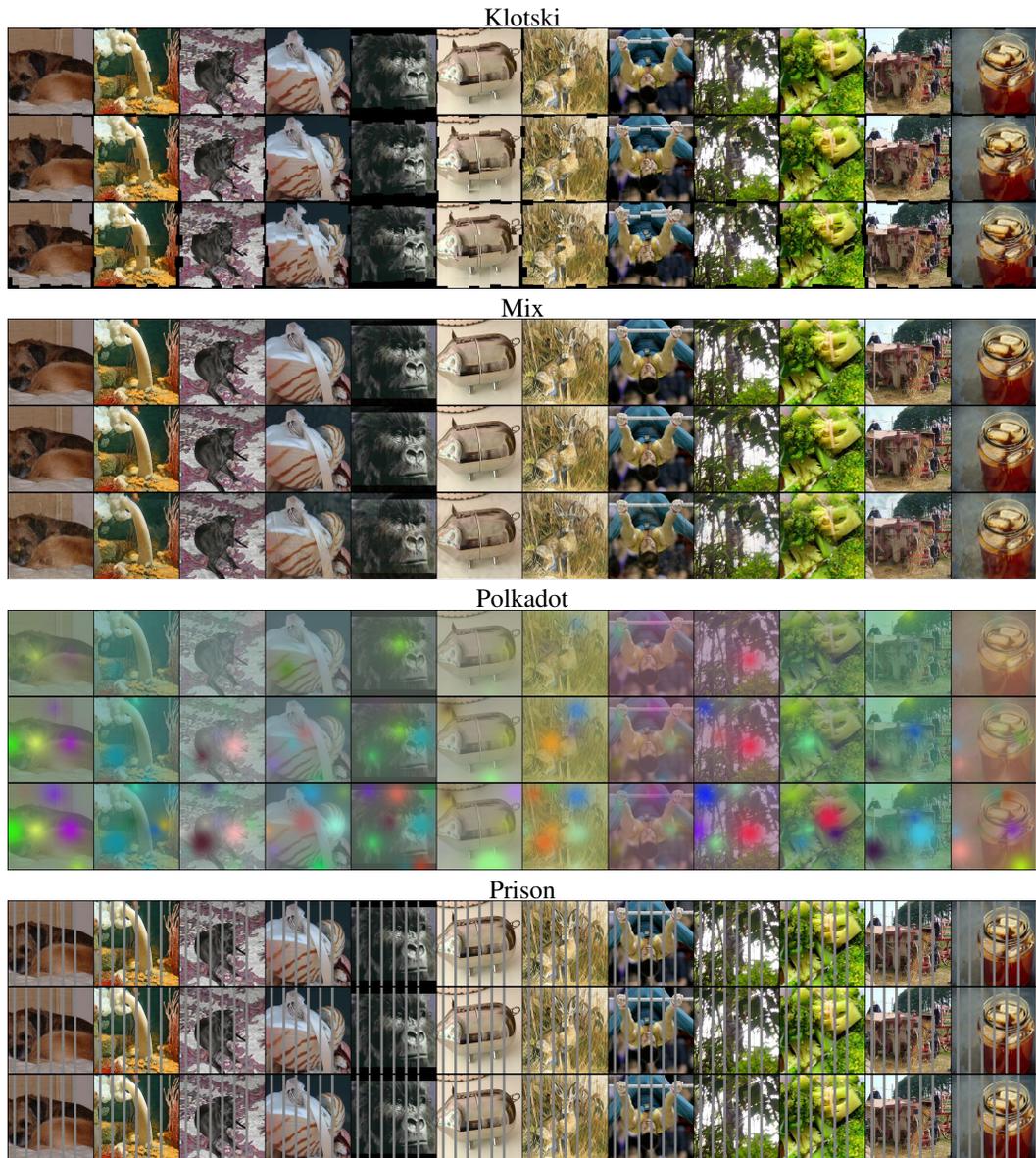


Figure 19: Attacked samples of low distortion (1st row), medium distortion (2nd row), and high distortion (last row) on a standard ResNet50 model

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

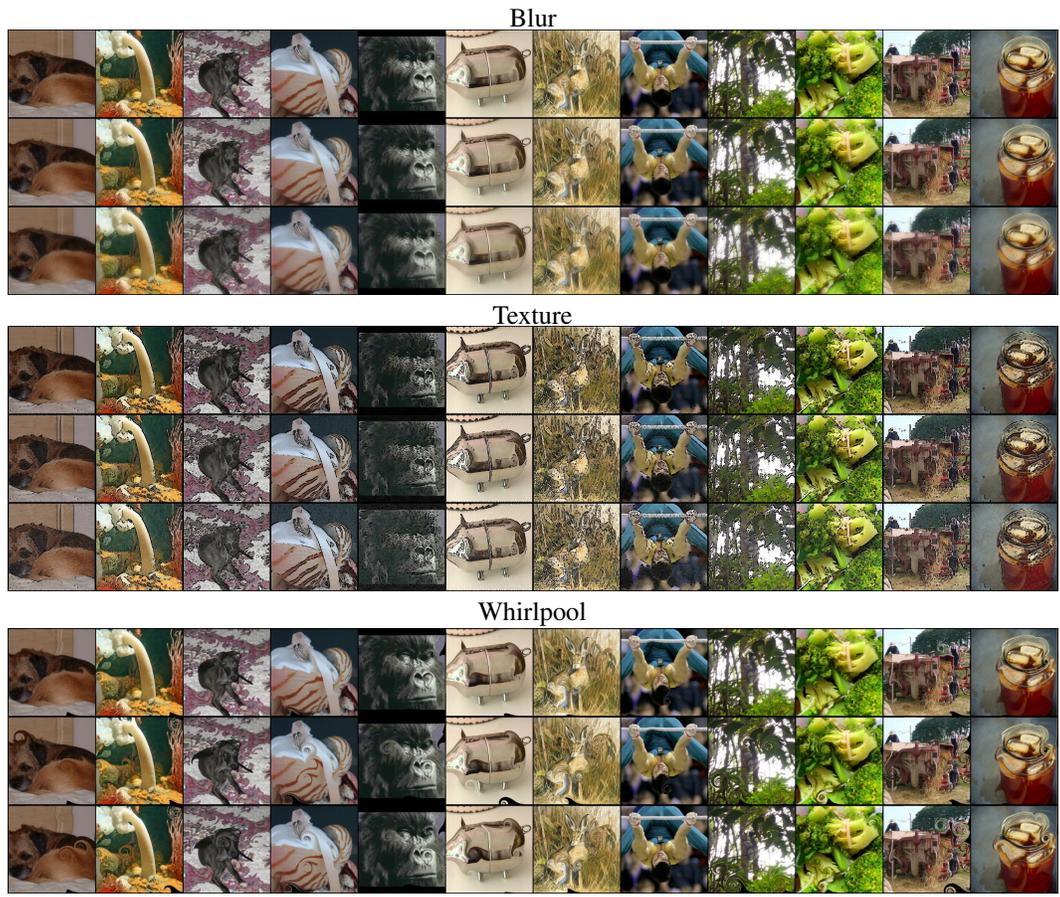


Figure 20: Attacked samples of low distortion (1st row), medium distortion (2nd row), and high distortion (last row) on a standard ResNet50 model

H SCALING BEHAVIOUR OF OUR ATTACKS

To see how our attacks perform across model scale, we make use of the ConvNeXt-V2 model suite (Woo et al., 2023) to test the performance of our attacks as we scale model size. We find that capacity improves performance across the board, but find diminishing returns to simply scaling up the architectures, pointing towards techniques described in Section 4.2.

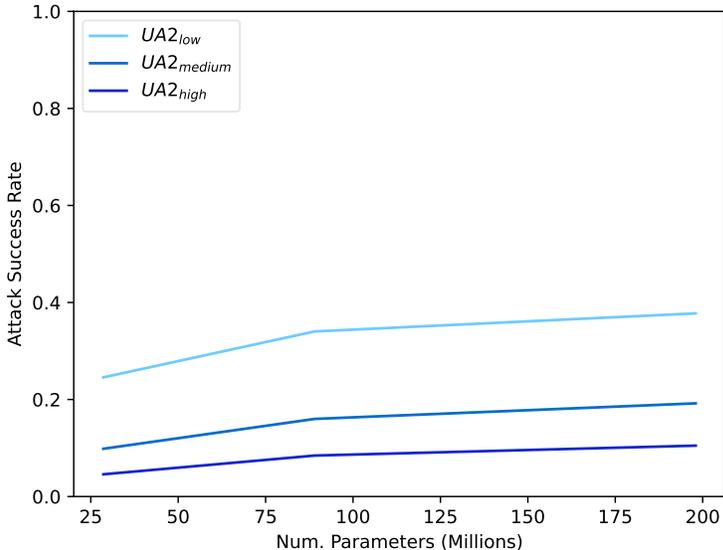


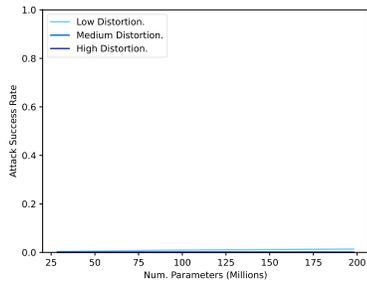
Figure 21: **Unforeseen Robustness across model scale.** We measure UA2 across model scale by evaluating the performance of ConvNeXt-V2 (Woo et al., 2023) models on ImageNet-UA, finding that scale improves performance, although the benchmark still provides a challenge to the largest models.

I DISTRIBUTION SHIFT COMPARED TO UNFORESEEN ROBUSTNESS

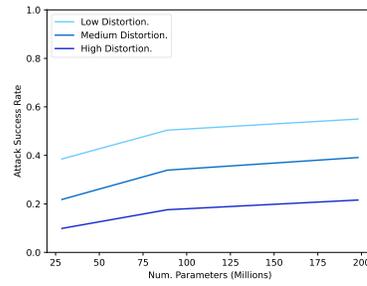
Table 13: **Distribution-shift benchmarks and UA2** Comparing performance on ImageNet-Sketch and ImageNet-R to performance against both non-optimized and optimized versions of UA2. We observe that performance on standard distribution shift benchmarks is correlated with performance on non-optimized UA2, while optimized UA2 settings favor models which have been trained for worst-case settings.

Model	UA2 (non-optimised)	ImageNet-Sketch Acc.	ImageNet-R Acc.	UA2
Resnet 50	55.2	24.1	36.2	1.6
Resnet50 + AugMix	59.1	28.5	41.0	3.5
Resnet50 + DeepAug	60.2	29.5	42.2	3.0
Resnet50 + Mixup	59.9	26.9	39.6	4.8
Resnet50 + L_2 , ($\epsilon = 5$)	43.2	24.2	38.9	13.9
Resnet50 + L_∞ , ($\epsilon = 8/255$)	40.6	18.6	34.8	10

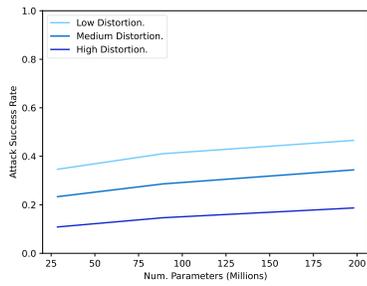
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673



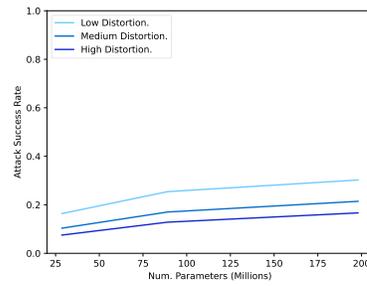
(a) JPEG



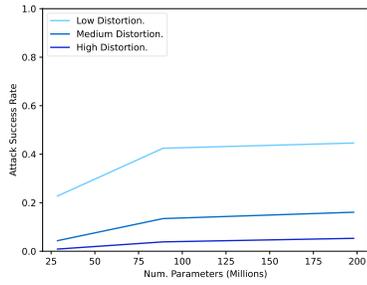
(b) Elastic



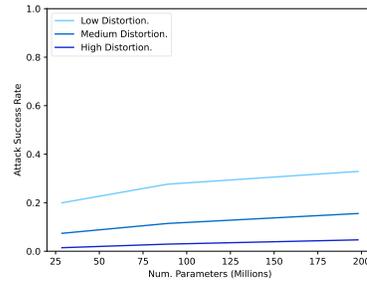
(c) Wood



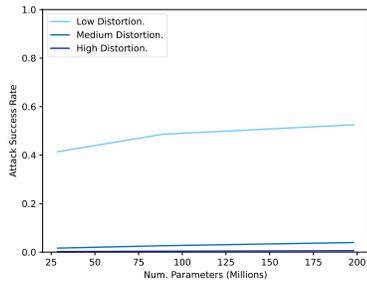
(d) Glitch



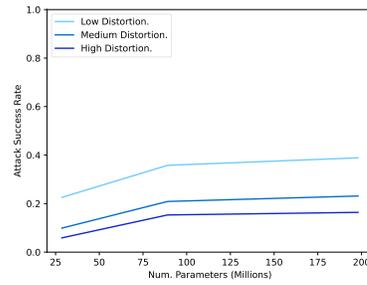
(e) Kaleidoscope



(f) Pixel



(g) Snow



(h) Gabor

Figure 22: **Behaviour of core attacks across model scale.** We see the performance of the eight core attacks across the ConvNeXt-V2 model suite, with performance on attacks improving with model scale.

J BENCHMARKING NON- L_p ADVERSARIAL TRAINING STRATEGIES

We wish to compare training strategies which have been specifically developed for robustness against both a variety of and unforeseen adversaries. To this end, we use Meta Noise Generation (Madaan et al., 2021b) as a strong multi-attack robustness baseline, finding that on CIFAR-10-UA this leads to large increases in robustness (Table 14). We also evaluate Perceptual Adversarial Training (Laidlaw et al., 2020) and Variational Regularization (Dai et al., 2022), two techniques specifically designed to achieve unforeseen robustness. We also evaluate combining PixMix and L_p adversarial training. All of these baselines beat L_p training.

Table 14: **Comparing alternative training strategies to L_p baselines** We demonstrate that models trained using Meta Noise Generation (MNG) (Madaan et al., 2021b) improve over L_p training baselines on CIFAR-10-UA.

Training	Clean Acc.	UA2
Standard	95.8	7.4
$L_\infty, \epsilon = 8/255$	86.5	39.8
$L_2, \epsilon = 2$	95.5	21.4
MNG	88.9	51.1

Meta Noise Generation (MNG) out-performs L_p baselines. We find that MNG, a technique original developed for multi-attack robustness shows a 11.3% increase in UA2 on CIFAR-10-UA, and PAT shows a 3.5% increase in UA2.

Table 15: **Specialised Unforeseen robustness training strategies.** We see that ImageNet-UA PAT (Laidlaw et al., 2020) and PAT-VR (Dai et al., 2022) trained ResNet50s improve over L_p baselines. Selected L_p models are the best Resnet50s from the bench-marking done in Figure 8, and for computational budget reasons they are trained on a 100-image subset of ImageNet, constructed by taking every 10th class.

Training	Clean Acc.	UA2
Standard	88.7	3.2
$L_\infty, \epsilon = 8/255$	79.7	17.5
$L_2, \epsilon = 4800/255$	71.6	25.0
PAT	75.0	26.2
PAT-VR	69.4	29.5

Table 16: **PixMix and L_p training.** We compare UA2 performance on CIFAR-10 of models trained with PixMix and adversarial training. Combining PixMix with adversarial training results in large improvements in UA2, demonstrating an exciting future direction for improving unforeseen robustness. All numbers denote percentages, and L_∞ training was performed with the TRADES algorithm.

Model	Clean Acc.	UA2
WRN-40-2 + PixMix	95.1	15.00
WRN-28-10 + L_∞ 4/255	89.3	37.3
WRN-28-10 + L_∞ 4/255 + PixMix	91.4	45.1
WRN-28-10 + L_∞ 8/255	84.3	41.4
WRN-28-10 + L_∞ 8/255 + PixMix	87.1	47.4

K HUMAN STUDY OF SEMANTIC PRESERVATION

Table 17: **Results of user study.** We run a user study on the 200 class subset of ImageNet presented as part of ImageNet-R (Hendrycks et al., 2021), assessing the multiple-choice classification accuracy of human raters, allowing raters to choose certain images as corrupted. We use 4 raters per label and take a majority vote, finding high classification accuracy across all attacks.

Attack Name	Correct	Corrupted or Ambiguous
Clean	95.4	4.2
Elastic	92.0	2.0
Gabor	93.4	4.0
Glitch	80.2	16.0
JPEG	93.4	0.6
Kaleidoscope	93.0	6.2
Pixel	92.6	1.8
Snow	90.0	3.2
Wood	91.4	1.8
Adversarial images average	91.2	4.5

We ran user studies to compare the difficulties of labeling the adversarial examples compared to the clean examples. We observe that under our distribution of adversaries users experience a 4.2% drop in the ability to classify. This highlights how overall humans are still able to classify over 90% of the images, implying that the attacks have not lost the semantic information, and hence that models still have room to grow before they match human-level performance on our benchmark.

In line with ethical review considerations, we include the following information about our human study:

- **How were participants recruited?** We made use of the surgehq.ai platform to recruit all participants.
- **How were the participants compensated?** Participants were paid at a rate of \$0.05 per label, with an average rating time of 4 seconds per image—ending at an average rate of roughly \$45 hour.
- **Were participants given the ability to opt out?** All submissions were voluntary.
- **Were participants told of the purpose of their work?** Participants were told that their work was being used to “validate machine learning model performance”.
- **Was any data or personal information collected from the participants?** No personal data was collected from the participants.
- **Was there any potential risks done to the participants?** Although some ImageNet classes are sometimes known to contain elicit or unwelcome content Prabhu (2019). Our 100-class subset of ImageNet purposefully excludes such classes, and as such participants were not subject to any undue risks or personal harms.

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

Adversarial Images Classification

This work is used to validate machine learning model performance and your participation is voluntary. You're free to stop the task at any point in time. You'll be shown an image. One of the labels is indeed present in the image please select the correct one. If you're unfamiliar with a label take a second to search for it on google images. Please let us know if this happens often.

The image may however be too corrupted in which case select that it is too corrupted. **Please avoid using corrupted label unless necessary.**

Thanks!



Select which label is present in the image or if the image is too corrupted.

- Granny Smith (type of eating apple)
- pretzel (type of pretzel)
- pufferfish (type of fish)
- saxophone (type of musical instrument)
- accordion (type of musical instrument)
- Image is too corrupted

Next preview

Figure 23: **Interface of participants.** We demonstrate the interface which was provided to the participants of the study, involving the selection of correct classes from our 100-class subset of ImageNet.

This work is used to validate machine learning model performance and your participation is voluntary. You're free to stop the task at any point in time.

You'll be shown an image. One of the labels is indeed present in the image please select the correct one. If you're unfamiliar with a label take a second to search for it on google images. Please let us know if this happens often.

The image may however be too corrupted in which case select that it is too corrupted. Please avoid using corrupted label unless necessary.

Thanks!

Figure 24: **Instructions given to the participants.** Above is a list of the instructions which were given to the participants in the human study.

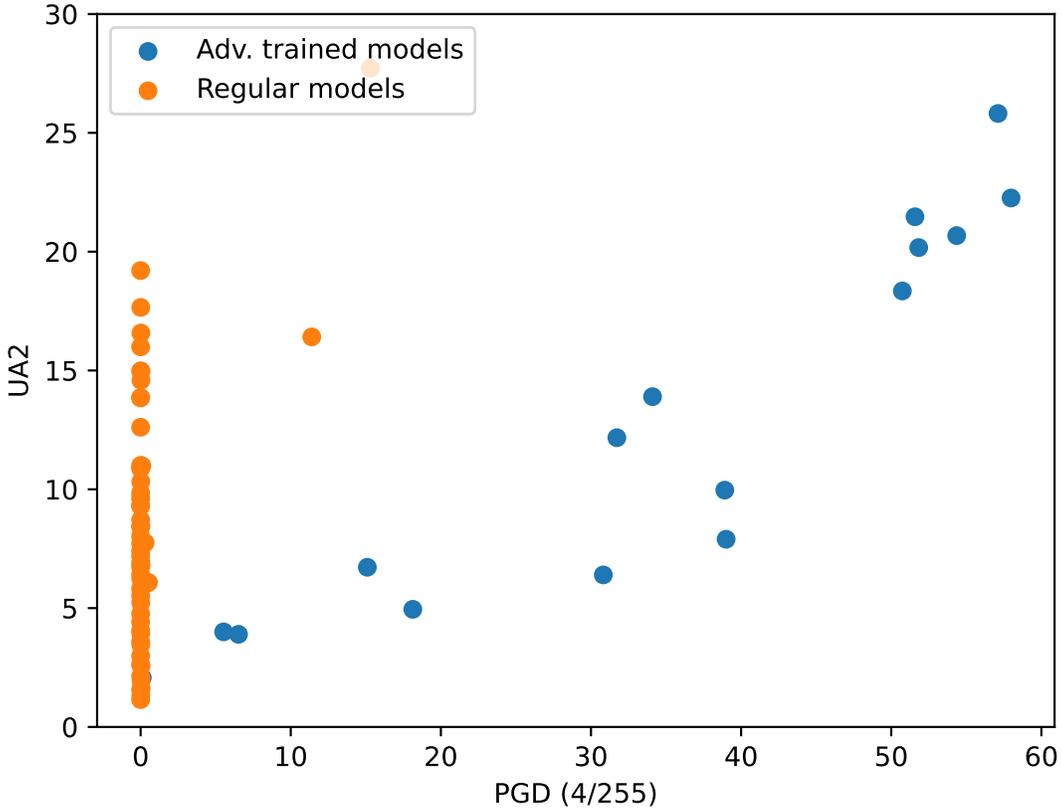
L CORRELATION OF L_p ROBUSTNESS AND ImageNet-UA

Figure 25: L_p **robustness correlates with UA2**. Across our benchmark, for adversarially trained models L_p robustness correlates with UA2 - however, several models trained without adversarial training still improve on UA2.

M GRID SEARCH VS. GRADIENT-BASED SEARCH

Table 18: **Comparing gradient-based search to grid-based search** We compare the performance of optimising with a randomised grid-based search using 1000 forward-passes per datapoint, finding that our gradient-based methods perform a lot better than this compute-intensive baseline.

Optimisation Technique	UA2
Randomized grid search	74.1
Gradient-based search (ours)	7.2

N TRANSFER ATTACKS

Table 19 shows the transfer-attack performances across various source and target models based on 1000 test samples. We observe that while the transfer attacks are not as effective as white-box attacks, they consistently outperform baseline unoptimized attacks where the perturbations are randomly initialized (Table 20).

Table 19: Transfer attack performance

	Clean Acc.	PGD	UA2	JPEG	Elastic	Wood	Glitch	Kal.	Pixel	Snow	Gabor
ResNet50 (source model)	75.2	0	13.2	0	22.2	30.8	10	4.3	4.8	3.1	30.4
ViT-small Patch16 ImageNet1K	78.5	73.1	59.99	75	62.7	69.9	46	48	62.8	55.5	60
ConvNeXt-V2-tiny ImageNet1K	82.1	74.8	67.66	77.1	69	75.9	54	60	73.6	65.2	66.5
Swin-small ImageNet1K + L_∞ 4/255	71.1	70.6	50.39	70.9	56.7	65.8	34.8	10.7	59.3	48.4	56.5
ResNet50	75.2	67.9	43.19	70.1	53.1	57.7	30.1	5.4	53.3	38.1	37.7
ViT-small Patch16 ImageNet1K (source model)	78.5	0	6.51	0	8.2	12.7	0.5	4.7	2.1	0.8	23.1
ConvNeXt-V2-tiny ImageNet1K	82.1	75.7	67.3	78.6	68.5	72.8	56.4	59.9	70.1	65.1	67
Swin-small ImageNet1K + L_∞ 4/255	71.1	70.5	50.11	70.9	57.1	65.1	35	10.8	59.5	48	54.5
ResNet50	75.2	67.8	42.06	68.3	51	55.7	31.7	5.8	51.7	32.1	40.2
ViT-small Patch16 ImageNet1K	78.5	74.7	57.31	75	60	69	42	46.8	57.2	50.2	58.3
ConvNeXt-V2-tiny ImageNet1K (source model)	82.1	0	12.15	0	23.2	22.3	7.4	3.5	6	0.6	34.2
Swin-small ImageNet1K + L_∞ 4/255	71.1	71.2	50.1	71.2	56.1	65	37.8	10.7	59.1	45	55.9
ResNet50	75.2	64	36.95	61.8	42.5	57.8	15.6	5.4	45.3	29.2	38
ViT-small Patch16 ImageNet1K	78.5	66.9	53.3	70.6	51.4	68.2	23.8	47.1	58.4	44.2	62.7
ConvNeXt-V2-tiny ImageNet1K	82.1	75.5	65.26	75.7	64.7	74.5	46.1	58.2	72.3	63.6	67
Swin-small ImageNet1K + L_∞ 4/255 (source model)	71.1	53.8	21.4	42	17.9	42.3	5.1	5.1	7.6	3.4	47.8

Table 20: Unoptimized attack performance

	Clean Acc.	PGD	UA2	JPEG	Elastic	Wood	Glitch	Kal.	Pixel	Snow	Gabor
ResNet50	75.2	74.1	56.44	74.3	62.8	55.7	55.8	6.3	74.1	74.8	47.7
ViT-small Patch16 ImageNet1K	78.5	78	69.19	78	70.2	70.2	65.4	47.7	77.3	78.6	66.1
ConvNeXt-V2-tiny ImageNet1K	82.1	82.2	74.74	82.2	75.2	74.4	69.7	60.7	81.5	81.4	72.8
Swin-small ImageNet1K + L_∞ 4/255	71.1	71.3	58.19	71.6	62	63.4	58	10.2	70.9	71.7	57.7

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

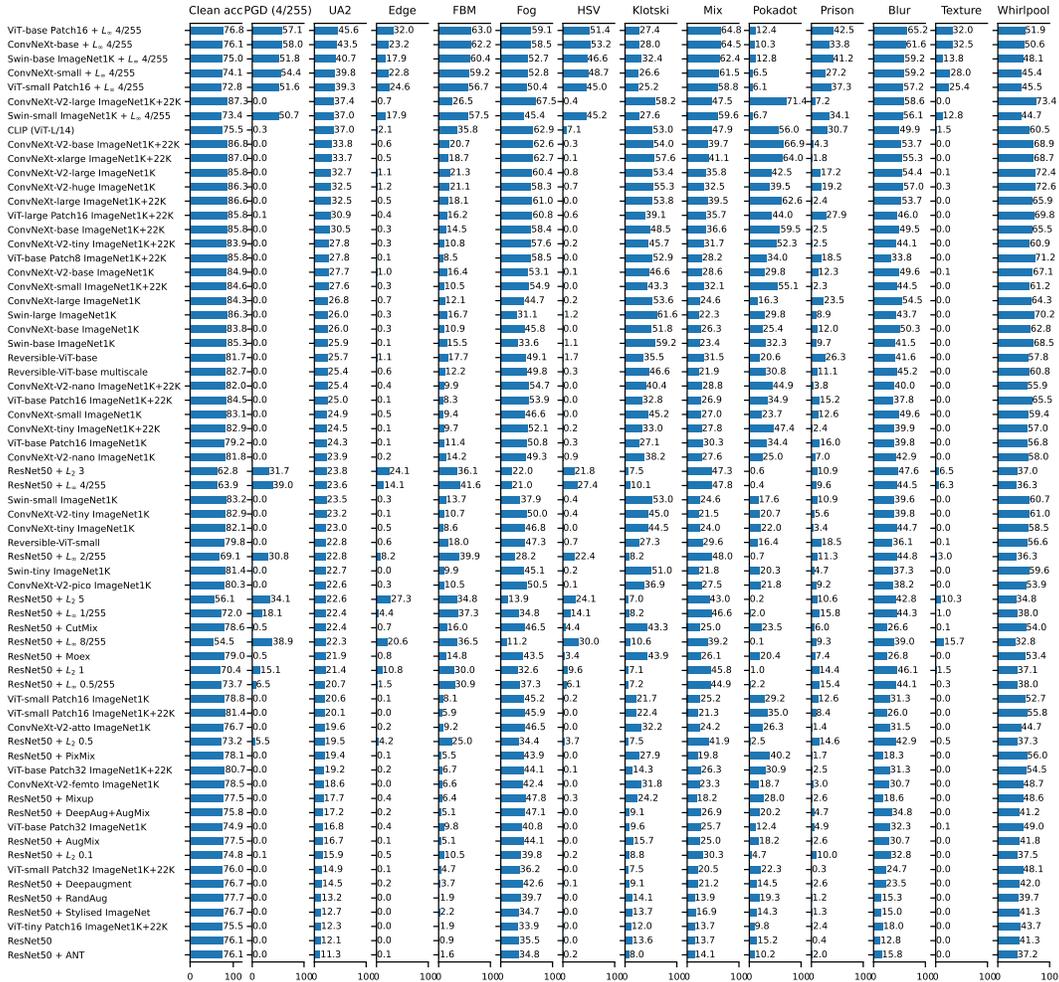


Figure 26: ImageNet UA2 performance under extra attacks in medium distortion