

GENERALIZABILITY OF ADVERSARIAL ROBUSTNESS UNDER DISTRIBUTION SHIFTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent progress in empirical and certified robustness promises to deliver reliable and deployable Deep Neural Networks (DNNs). Despite that success, most existing evaluations of DNN robustness have been done on images sampled from the same distribution that the model was trained on. Yet, in the real world, DNNs may be deployed in dynamic environments that exhibit significant distribution shifts. In this work, we take a first step towards thoroughly investigating the interplay between empirical and certified adversarial robustness on one hand and domain generalization on another. To do so, we train robust models on multiple domains and evaluate their accuracy and robustness on an unseen domain. We observe that: (1) both empirical and certified robustness generalize to unseen domains, and (2) the level of generalizability does not correlate well with input visual similarity, measured by the FID between source and target domains. We also extend our study to cover a real-world medical application, in which adversarial augmentation enhances both the robustness and generalization accuracy in unseen domains.

1 INTRODUCTION

Deep Neural Networks (DNNs) are vulnerable to small and carefully designed perturbations, known as adversarial attacks (Szegedy et al., 2014; Goodfellow et al., 2015). That is, a DNN $f_\theta : \mathbb{R}^d \rightarrow \mathcal{P}(\mathcal{Y})$ can produce two different predictions for the inputs $x \in \mathbb{R}^d$ and $x + \delta$, although both x and $x + \delta$ are perceptually indistinguishable. Furthermore, DNNs are found to be brittle against simple semantic transformations such as rotation, translation, and scaling (Engstrom et al., 2019). These observations raised concerns regarding the deployability of DNNs in security-critical applications, such as self-driving and medical diagnosis (Papernot et al., 2016; Finlayson et al., 2019; Ma et al., 2021). This brittleness provoked several efforts to build models that are not only accurate but also *robust* (Gu & Rigazio, 2015). Building robust models is usually achieved either (i) *empirically*, where the DNN’s training routine is changed to include such malicious adversarial examples in the training set (Madry et al., 2018), or (ii) *certifiably*, where theoretical guarantees are given about the robustness of a DNN in a region around a given input x (Lécuyer et al., 2019).

Despite great progress in the adversarial robustness literature on building accurate and robust models, most approaches are tested on *in-distribution* data. In other words, the scenario considered is one in which both the training and testing sets are independently and identically distributed (IID). However, this IID assumption rarely holds in practice, as data in the real world can be sampled from various distributions with significant domain shifts. For example, a deep-learning based medical image classifier may be trained on data collected from one hospital, but later deployed in a different hospital (Bánci et al., 2019). Unfortunately, DNNs struggle to generalize to out-of-domain data (Geirhos et al., 2020; 2021), even in the absence of adversarial examples. This lack of generalization has led the research community to invest in the problem of Domain Generalization (DG). The aim of DG is to learn invariant representations from diverse distributions of data, denoted as *source* domains, such that these representations generalize to an unseen distribution, known as the *target* domain (Wang et al., 2021; Gulrajani & Lopez-Paz, 2021). This setup mimics the unexpected nature of real-world distribution shifts, where models are constantly exposed to novel domains, and fine-tuning on all these domains becomes impractical. While there has been considerable effort in improving the generalization of DNNs (Tzeng et al., 2014; Sun & Saenko, 2016; Motiian et al., 2017; Zhang et al., 2021; Shen et al., 2021; Wang et al., 2021; Zhou et al., 2022), the generalizability of adversarial robustness to unseen domains remains unexplored.

In this work, we set out to study the interplay between domain generalization and adversarial robustness. We conduct comprehensive experimental studies on five standard DG benchmarks provided by DomainBed (Gulrajani & Lopez-Paz, 2021) and WILDS (Koh et al., 2021). In our experiments, we study both empirical and certified robustness against input perturbations and spatial deformations. We first investigate the generalizability of empirical robustness, which a DNN obtains by employing the popular adversarial training method (Madry et al., 2018) while training on the source data. We find that, in many scenarios, improving empirical robustness in the source domain generalizes to the target domain with little cost on the performance of the model on unperturbed data. We then inspect the generalizability of certified robustness against both input perturbations and parametric deformations by employing Randomized Smoothing (RS) (Cohen et al., 2019) and DeformRS (Alfarra et al., 2022a). We observe that certified robustness generalizes to unseen domains when using randomized smoothing frameworks against pixel perturbations and five different input deformations including rotation, translation, and scaling. To the best of our knowledge, we provide the first large scale experimental analysis of the generalizability of adversarial robustness to unseen domains. Our analysis leads to the following contributions:

1. We contrast the behavior of robustness under both transfer learning and domain generalization. Unlike transfer learning, domain generalization does not always improve through robust training.
2. We empirically show that visual similarity, between the source and target domains, does not correlate well with the level of generalizability to the target domain.
3. We analyze a practical medical application, in which adversarial training in the source domain improves the generalization of accuracy and robustness in the target domain.

2 RELATED WORK

Domain Generalization. Domain generalization (DG) studies the ability of models to learn representations that can be readily applied to data from unseen domains (Wang et al., 2021; Zhou et al., 2022). In the DG setup, a model is trained on multiple source domains and then evaluated on an unseen target domain, which exhibits a significant shift from the training domains. DG approaches can be categorized into different groups. (i) Data augmentation techniques learn generalizable models by increasing the diversity of the source data (Gong et al., 2019; Zhou et al., 2020; 2021). (ii) Representation learning methods aim at extracting domain-invariant representations that seamlessly apply in any unseen domain (Blanchard et al., 2011; Nguyen et al., 2021; Lu et al., 2022) (iii) Learning-strategy approaches may achieve generalization through meta-learning, self-supervised learning, or optimization procedures that seek flat minima (Li et al., 2018; Carlucci et al., 2019; Cha et al., 2021). In this work, we study DG from an adversarial robustness lens. In particular, we analyze both the generalization accuracy and robustness of adversarially trained classifiers on unseen domains.

Adversarial Robustness. Adversarial attacks are imperceptible, semantic-preserving perturbations that can fool DNNs (Goodfellow et al., 2015; Szegedy et al., 2014). Given the security concerns that adversarial attacks induced, several works proposed changing the training routine to enhance model robustness (Zhang et al., 2019). For example, adversarial training (Madry et al., 2018) encourages the model to classify adversarial examples correctly. While empirical defenses like adversarial training are effective in enhancing the robustness of the underlying model, such approaches do not guarantee robustness. Subsequently, many empirical defenses were broken when more powerful attacks were designed (Carlini & Wagner, 2017; Athalye et al., 2018). As a result, there has been a growing interest in certifiably robust classifiers, for which no adversary can exist in a specified region around a data point (Raghunathan et al., 2018; Mohapatra et al., 2020; Lee et al., 2021). A scalable approach to achieving certified robustness is Randomized Smoothing (RS) (Cohen et al., 2019). RS constructs a smooth classifier from any arbitrary base classifier by outputting the most probable class when the input is subjected to Gaussian noise. Recently, DeformRS extended RS to provide certified robustness against parameterized geometric deformations (Alfarra et al., 2022a; S. et al., 2022). In this work, we set out to study the interplay between (empirical and certified) robustness and domain generalization by deploying adversarial training, RS, and DeformRS.

Adversarial Training in Dynamic Environments. To improve the ability of machine learning models to learn generalizable knowledge, researchers have proposed several problems, such as trans-

Problem Setup	Training Data	Target Data	Problem Condition	Access to Target
Transfer learning	S_{source}, S_{target}	S_{target}	$\mathcal{Y}_{source} \neq \mathcal{Y}_{target}$	✓
Domain generalization	$S = \{S_i i = 1, \dots, N\}$	S_{N+1}	$\mathbb{P}_{XY}(S_k) \neq \mathbb{P}_{XY}(S_n)$ for $k \neq n$	✗

Table 1: **Comparison between Domain Generalization (DG) and Transfer Learning.** DG differs from transfer learning in two main ways. 1) The model in DG never sees the target data during training, so fine-tuning on the target is not allowed. 2) The target labels are kept fixed in DG; however, the target samples are drawn from a domain that is distinct from the source domains.

fer learning, continual learning, domain adaptation, and domain generalization (Zhuang et al., 2020; Delange et al., 2021; Wang & Deng, 2018; Wang et al., 2021). Among these problems, only transfer learning, where a model pre-trains on tasks with large datasets and then adapts to downstream tasks with limited data, has been thoroughly studied under the lens of adversarial robustness. Salman et al. (2020) showed that, in terms of downstream task accuracy, adversarially trained representations outperform nominally trained representations. Utrera et al. (2021) further explained that adversarial training in the source domain increases shape bias, resulting in better transferability. Finally, Deng et al. (2021) provided theoretical justification to support these empirical findings. Besides downstream task accuracy, Shafahi et al. (2020) studied the transferability of robustness itself. Although useful, these transfer learning results presume fine-tuning on the target domain, which is not possible in many real-life scenarios. Table 1 illustrates the differences between transfer learning and domain generalization, which is the setup we adopt. In this paper, we take a first step to empirically investigate whether adversarial training leads to robust representations that generalize well without requiring prior knowledge of the target domain.

3 BACKGROUND ON DOMAIN GENERALIZATION

Domain Generalization Setup. Given an input space \mathcal{X} and a label space \mathcal{Y} , one can define a joint distribution \mathbb{P}_{XY} over \mathcal{X} and \mathcal{Y} . A domain, or distribution, is a collection of samples drawn from \mathbb{P}_{XY} . In our setting, the input and label spaces are fixed, but we may have multiple unique joint distributions. Specifically, we assume that there are N source domains of varying sizes. For each task n , $S_n = \{(x_j, y_j)\}_{j=1}^{|S_n|} \sim \mathbb{P}_{XY}^{(n)}$. We denote the training set by $S = \{S_i | i = 1, \dots, N\}$. The aim of DG is to use S to learn a mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$ that minimizes the error on some unseen target domain $T \sim \mathbb{P}_{XY}^{(N+1)}$. We enforce that $\mathbb{P}_{XY}^{(k)} \neq \mathbb{P}_{XY}^{(n)}$ for $k \neq n, k, n \in \{1, \dots, N+1\}$, which means that the target domain is distinct from the source domains that are, in turn, also distinct from each other. More formally, given S , we seek a parameterized model f_{θ^*} such that:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathbb{P}_{XY}^{(T)}} [\mathcal{L}(f_{\theta}(x), y)], \quad (1)$$

where \mathcal{L} is the cross-entropy loss for the classification task. Note that the model is not allowed to sample the target domain during training, so most methods use the empirical risk of the source datasets as a proxy for the true target risk. The supervised average risk (\mathcal{E}) is given by:

$$\mathcal{E} = \frac{1}{N} \sum_{n=1}^N \frac{1}{|S_n|} \sum_{i=1}^{|S_n|} [\mathcal{L}(f_{\theta}(x), y)] \quad (2)$$

with $(x, y) \sim S$. In practice, we define a fixed held-out validation set $S^v \subset S$. The average risk on this source validation set is used to select the best model, which is evaluated on the target domain without any fine-tuning steps. In what follows, Section 4 (and Section 5) investigates the generalizability of empirical (and certified) robustness to diverse target domains.

4 EMPIRICAL ROBUSTNESS AND DOMAIN GENERALIZATION

In this section, we study the generalizability of empirical robustness methods that enhance the adversarial robustness of DNNs. We begin with a brief introduction of Adversarial Training (AT) (Madry et al., 2018), after which we study the effect of deploying AT in a domain generalization setup.

4.1 BACKGROUND AND SETUP

Adversarial Attacks. Adversarial attacks are small imperceptible perturbations that, once added to a “clean” input sample, cause the classifier f_θ to misclassify the perturbed sample. Formally, let (x, y) be an input label pair where f_θ correctly classifies x (i.e. $\arg \max_i f_\theta^i(x) = y$). An attacker crafts a small perturbation δ such that $\arg \max_i f_\theta^i(x + \delta) \neq y$, which is usually obtained by solving the following optimization problem:

$$\max_{\delta} \mathcal{L}(f_\theta(x + \delta), y) \quad \text{s.t. } \|\delta\|_p \leq \epsilon, \quad (3)$$

where $p \in \{2, \infty\}$, $\epsilon > 0$ is a small constant that enforces the imperceptibility of the added perturbation, and \mathcal{L} is a suitable loss function (e.g. Cross Entropy). Let δ^* be the solution to the problem in Eq. 3, then the adversarial example is denoted by $x_{adv} = x + \delta^*$.

Adversarial Training as Augmentation. Adversarial Training (AT) (Madry et al., 2018) trains the classifier on adversarial examples rather than clean samples. In particular, AT obtains the network parameters θ^* by solving the following optimization problem:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta, \|\delta\|_p \leq \epsilon} \mathcal{L}(f_\theta(x + \delta), y) \right], \quad (4)$$

where \mathcal{D} is a data distribution. In general, the inner maximization problem is solved through K steps of Projected Gradient Descent (PGD) (Madry et al., 2018). While conducting adversarial training enhances the model’s robustness against adversarial attacks, this usually comes at the cost of losing some clean accuracy (performance on unperturbed samples). To alleviate the drop in performance, we follow the method by Zhang et al. (2019) and deploy adversarial training as a data augmentation scheme. In particular, we obtain network parameters θ^* that minimize the following objective:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{(x,y) \sim \mathbb{P}_{XY}^{(T)}} [\lambda \mathcal{L}(f_\theta(x), y) + (1 - \lambda) \mathcal{L}(f_\theta(x_{adv}), y)], \quad (5)$$

where $\lambda \in [0, 1]$ controls the robustness-accuracy trade-off.

Training & Evaluation Setup. In our experiments, we focus on image classification and adopt the framework of DomainBed (Gulrajani & Lopez-Paz, 2021), which is the standard benchmark in the image domain generalization literature. All models are initialized with a ResNet-50 backbone pre-trained on ImageNet-1K (Deng et al., 2009). We train all models with adversarial augmentation to minimize the objective in Eq. 5 on the source domains, where x_{adv} is computed with a Projected Gradient Descent (PGD) attack (Madry et al., 2018) using 20 PGD steps. The target domain remains unseen until test time. Specifically, we follow the training-domain validation strategy described in DomainBed for model selection. We experiment with a range of perturbation budgets (ϵ) on various datasets: PACS, OfficeHome, VLCS, and TerraIncognita (Gulrajani & Lopez-Paz, 2021). We report the ℓ_∞ results in the main paper, where we use $\epsilon \in \{0, 8/255, 16/255, 32/255\}$. The appendix includes experiments for ℓ_2 perturbations due to space constraints. Note that for $\epsilon = 0$, the training objective reduces to the empirical risk minimization in Eq. equation 2. Each model is trained on one value of ϵ but is evaluated on all four values of ϵ with 20 steps of PGD attacks. In the following experiments, we fix λ in Eq. 5 to 0.5 and leave the ablation to the appendix.

4.2 THE GENERALIZATION OF EMPIRICAL ROBUSTNESS

In this section, we investigate the generalizability of empirical robustness to unseen domains. More precisely, we are interested in understanding the interplay between standard accuracy and robust accuracy in the scope of source vs. target domains. We report in Figure 1 the standard accuracy (first column of each matrix) and robust accuracy against different values of ϵ on all considered datasets. We further summarize the clean and robust accuracy at $\epsilon = 8/255$ in Table 2 for ease of comparison. Next, we analyze these results to answer the following questions.

Q1: Do adversarially robust models generalize better than their standard-trained counterparts? No, which is evident from the first two rows of Table 2. With the exception of the smallest dataset PACS, the clean accuracy of the robust model in the unseen domain is lower than that of the standard-trained model by 3.6% or more. **❶ Unlike transfer learning, where robust training in**

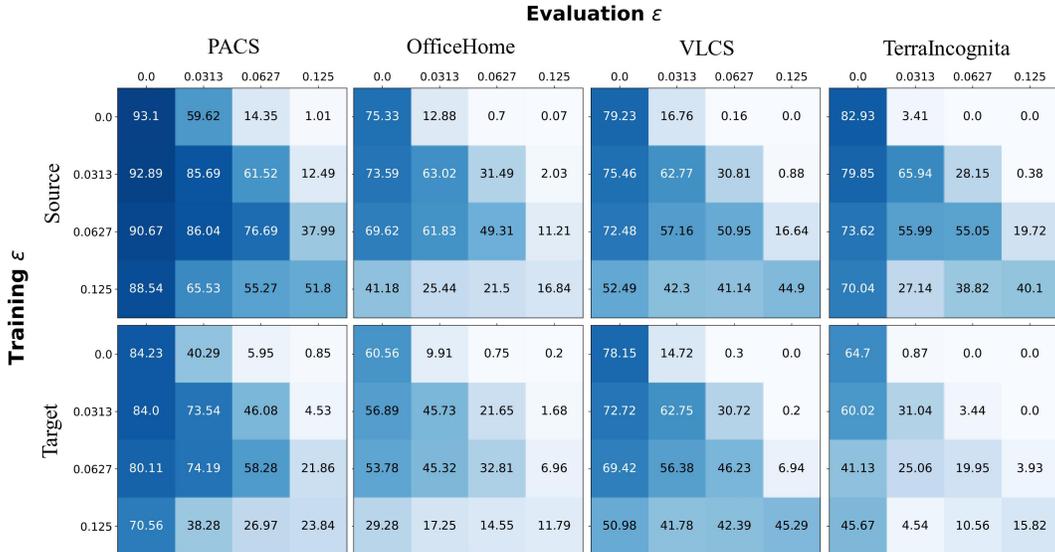


Figure 1: **Evaluation of ℓ_∞ Robustness.** We train models in the source domain and evaluate them in both the source and target domains with different ϵ values. In the Source tables, the $(i^{\text{th}}, j^{\text{th}})$ entry represents a model trained with $\epsilon = i$ and evaluated on $\epsilon = j$ on samples drawn from source domain. In the Target tables, the $(i^{\text{th}}, j^{\text{th}})$ entry is a model trained in the source with $\epsilon = i$ and evaluated on $\epsilon = j$ on samples drawn from target domain. No fine-tuning is done in the target domain.

		Dataset			
		PACS	OfficeHome	VLCS	TerraIncognita
Clean Accuracy (%)	Standard Model	84.23	60.56	78.15	64.70
	Robust Model	84.00	56.89	72.72	60.02
Robust Accuracy (%)	Standard Model	40.29	9.91	14.72	0.87
	Robust Model	73.54	45.73	62.75	31.04

Table 2: **Generalization of Robustly- and Nominally- Trained Models on Various Datasets.** Applying adversarial training on the source domain leads to significant improvements in the model’s robustness in the target domain, relative to the observed drop in the standard accuracy.

the source domain is favorable, robust training does not improve generalization to the target domain if no fine-tuning is allowed. This result contrasts with findings from the transfer learning literature, where models trained robustly in the source domain outperform standard-trained models across a variety of downstream tasks (Salman et al., 2020). It is especially surprising given that previous works suggest that robust training encourages shape bias over texture bias, hinting at better generalization (Geirhos et al., 2019; Utrera et al., 2021). Moreover, Deng et al. (2021) showed that adversarial training in the source domain results in provably better representations for fine-tuning on the target domain. Such seemingly contradictory findings can be reconciled by considering the key differences between transfer learning and domain generalization summarized in Table 1. Specifically, previous works in transfer learning assume that the model can sample the target domain at some point to perform fine-tuning. Since domain generalization does not allow access to the target domain, such benefits are not guaranteed. *We encourage future works to investigate under what conditions adversarial training helps the generalization accuracy with no fine-tuning on the target.*

Q2: Does a higher source domain robustness correspond to a higher target domain robustness? As expected, DNNs lose some robustness when evaluated on a target domain that is distinct from the training domains. This observation is evident by comparing any cell in the top row (Source) tables in Figure 1 with the corresponding cell in the second row (Target) tables. For example, the TerraIncognita model, which is trained and evaluated on $(\epsilon = 8/255)$ adversaries, loses around 35% accuracy when the distribution shifts to the target domain. Nevertheless, by observing that all the source and target tables have similar color trends, we find that **higher robustness in the source domain corresponds to higher robustness in the target domain.** Our results suggest that one way

to increase the out-of-distribution robustness of a deployed model is to improve its robustness in the source validation set, which supports the applicability of ongoing efforts in adversarial robustness research (Zhang et al., 2019; Wang et al., 2020; Wu et al., 2020).

Q3: Does the robustness-accuracy trade-off generalize to unseen domains? As observed in Figure 1, **achieving a more robust model comes at the cost of standard accuracy not only in the source domain, but also in the target domain.** Looking at OfficeHome, the target robust accuracy of a robust model (trained and evaluated on $\epsilon = 16/255$) is 50% more than that of the standard-trained model. Yet, the clean accuracy of the robust model is about 6% less than the standard model accuracy. In general, as the network becomes more robust to a particular perturbation budget ϵ in the source domain, it becomes more robust to adversaries generated within that budget in the target domain. Nevertheless, the performance of the robust network on clean samples decreases in both domains. Therefore, *consistent with the robustness literature (Tsipras et al., 2019), robustness comes at the cost of standard accuracy even in the unseen target domains.*

5 CERTIFIED ROBUSTNESS AND DOMAIN GENERALIZATION

In Section 4, we analyzed the interplay between empirical robustness (obtained by adversarial training) and domain generalization. While empirical robustness studies give hints about the reliability of a given model under adversarial attacks, they give no guarantees against the existence of such adversaries. To deploy DNNs in dynamic environments (Koh et al., 2021), we need robustness guarantees to carry over into unseen domains. To that end, we study the generalizability of the certified robustness of DNNs. We deploy Randomized Smoothing (RS) and DeformRS to certify DNNs against input perturbations and deformations. We start by giving a brief overview of RS and DeformRS.

5.1 BACKGROUND AND SETUP

Certifying Against Additive Perturbations and Input Deformations. Randomized smoothing (RS) (Cohen et al., 2019) is a method for constructing a “smooth” classifier from a given classifier f_θ . The smooth classifier returns the average prediction of f_θ when the input x is subjected to additive Gaussian noise:

$$g_\theta(x) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma^2 I)} [f_\theta(x + \epsilon)]. \quad (6)$$

Let g_θ predict label c_A for input x with some confidence, *i.e.* $\mathbb{E}_\epsilon [f_\theta^{c_A}(x + \epsilon)] = p_A \geq p_B = \max_{c \neq c_A} \mathbb{E}_\epsilon [f_\theta^c(x + \epsilon)]$, then, as shown by Zhai et al. (2020), g_θ ’s prediction is certifiably robust at x with certification radius:

$$R = \frac{\sigma}{2} (\Phi^{-1}(p_A) - \Phi^{-1}(p_B)), \quad (7)$$

where Φ^{-1} is the inverse CDF of the standard Gaussian distribution. As a result of Eq. 7, $\arg \max_i g_\theta^i(x + \delta) = \arg \max_i g_\theta^i(x)$, $\forall \|\delta\|_2 \leq R$.

While Eq. 7 provides theoretical guarantees for robustness against additive perturbations, DNNs are also brittle against simple input transformations such as rotation. Alfarrar et al. (2022a) extended randomized smoothing to certify parametric input deformations through DeformRS, which defined the parametric smooth classifier for a given input x with pixel coordinates p as follows:

$$g_\phi(x, p) = \mathbb{E}_{\epsilon \sim \mathcal{D}} [f_\theta(I_T(x, p + \nu_{\phi+\epsilon}))], \quad (8)$$

where I_T is an interpolation function (*e.g.* bilinear interpolation) and ν_ϕ is a parametric deformation function with parameters ϕ (*e.g.* ν is a rotation function and ϕ is the rotation angle). Analogous to the RS formulation in Eq. 6, g outputs the average prediction of f_θ over deformed versions of the input x . Alfarrar et al. (2022a) showed that parametric-domain smooth classifiers are certifiably robust against perturbations to the parameters of the deformation function. In particular, g ’s prediction is constant with certification radius:

$$\begin{aligned} R &= \sigma (p_A - p_B) && \text{for } \mathcal{D} = \mathcal{U}[-\sigma, \sigma], \\ R &= \frac{\sigma}{2} (\Phi^{-1}(p_A) - \Phi^{-1}(p_B)) && \text{for } \mathcal{D} = \mathcal{N}(0, \sigma^2 I). \end{aligned} \quad (9)$$

Put simply, as long as the perturbations to the deformation function parameters (*e.g.* rotation angle) are within R , the prediction of g remains constant. In this work, we leverage RS and DeformRS to study the generalizability of certified robustness to unseen target domains.

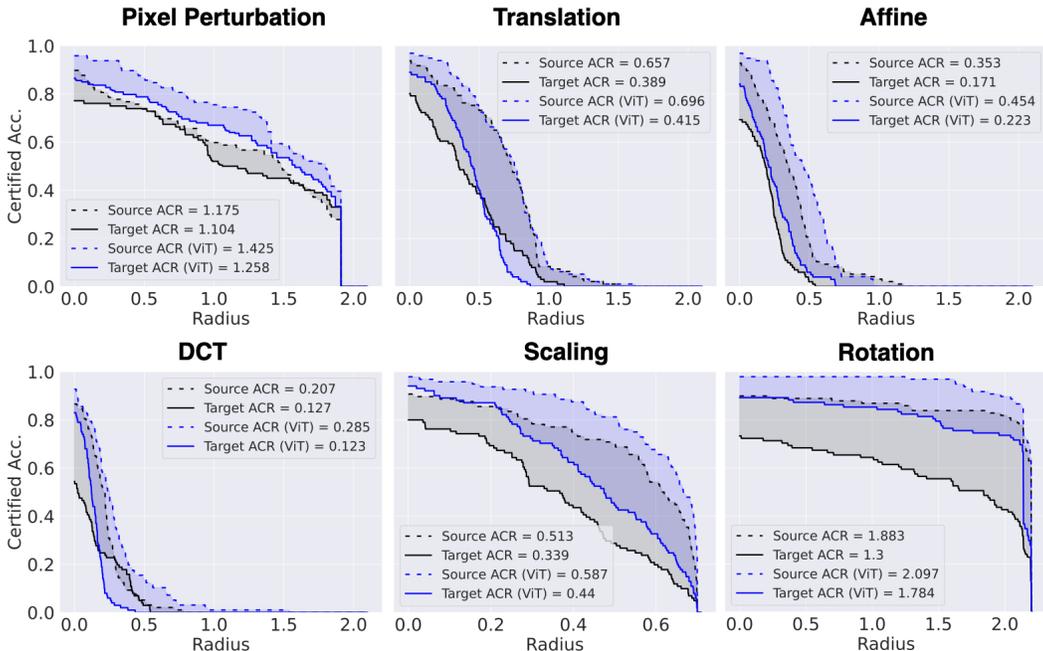


Figure 2: **Generalizability of Certified Robustness.** We certify ResNet-50 and ViT-Base against pixel perturbations and input deformations in the source and target domains of PACS. We observe that 1) certified robustness deployed via randomized smoothing generalizes to unseen domains, and that 2) a stronger architecture (ViT-Base) leads to a better source and target certified accuracy.

Experimental Setup. To split the data into source and target domains, we use the *Photo*, *Art*, *Cartoon*, and *Sketch* distributions from PACS (Li et al., 2017). We use RS to certify pixel perturbations and DeformRS to certify five input deformations: rotation, translation, scaling, affine, and a deformation characterized by a Discrete Cosine Transform (DCT) basis. Following (Gulrajani & Lopez-Paz, 2021), we employ data augmentation during training and train solely on the source domains. To evaluate the certified robustness of the trained classifier, we plot the certified accuracy curves for both the source and target domains for each considered deformation. The certified accuracy at a radius R is the percentage of the test set that is both classified correctly and has a certified radius of at least R . We calculate the certified radius for a given input through either Eq. 7 for pixel perturbations or Eq. 9 for input deformations. Here, we report the envelope plots, which illustrate the best certified accuracy per radius over all values of the smoothing deformation parameter ϕ . We leave the detailed results for each choice of ϕ to the appendix. We employ Monte Carlo sampling with 100k samples to estimate p_A and bound $p_B = 1 - p_A$ by following the standard practice (Zhai et al., 2020; Cohen et al., 2019; Alfara et al., 2022a). Finally, we follow (Zhai et al., 2020) in reporting the Average Certified Radius (ACR) of correctly classified samples.

Regarding the architecture, we follow the DomainBed (Gulrajani & Lopez-Paz, 2021) benchmark in selecting ResNet-50 as a backbone. To assess the effect of deploying a more powerful architecture on the generalizability to unseen domains, we further include experiments with the recent transformer model ViT-Base (Dosovitskiy et al., 2021).

5.2 GENERALIZABILITY OF CERTIFIED ROBUSTNESS TO UNSEEN TARGET DOMAINS

We investigate under what scenarios the certified robustness generalizes to unseen domains. We first show how much certified accuracy (CA) is maintained when the target domain exhibits a distribution shift. Then, we study whether a stronger backbone architecture can boost the CA generalizability. Finally, we evaluate how well perceptual similarity, as measured by FID and R-FID (Heusel et al., 2017; Alfara et al., 2022c), predicts the generalization of certified robustness.

Q4: Can certified robustness, obtained via randomized smoothing, generalize to unseen domains? We train smooth classifiers on a collection of source domains and certify the models on

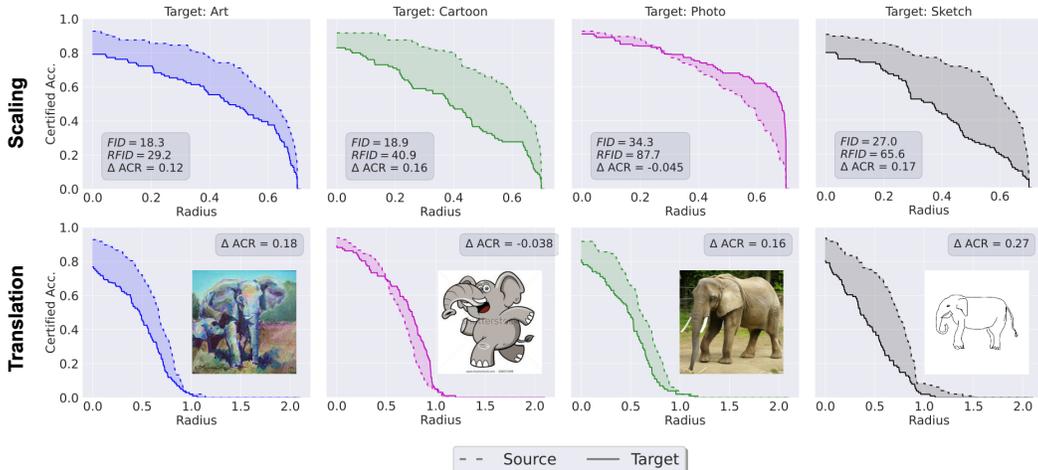


Figure 3: **Does visual similarity correlate with robustness generalizability?** We vary the target domain and plot the certified accuracy curves for two deformations: scaling and translation. A sample from each domain is shown in the second row. The FID/R-FID distances between the source domains and each target are reported in the first row. Visual similarity, measured by FID and R-FID, does not correlate with the level of robustness generalization to the target domain.

both the source and target domains. The target domains are *unseen* before certification. We plot the source CA curve with dashed black lines and the target CA curve with solid blue in Figure 2, along with the corresponding ACR. Our results show that **4 a considerable portion of the certified robustness, acquired by randomized smoothing, is maintained in the unseen domain.** When certified against pixel perturbations in the unseen domain, the average certified radius of ResNet-50 drops by around 6% only. Utilizing DeformRS, we extend this result from simple pixel perturbations to geometric deformations, like rotation and translation. This experiment is promising, since the models are never trained on the target data, but still exhibit some certified robustness. This validates the importance of recent research efforts that improve on randomized smoothing (Zhai et al., 2020; Alfarrar et al., 2022b). *To address real-world security challenges, we encourage future certified robustness works to conduct experiments on domain generalization datasets.*

Q5: Does the target certified accuracy improve when the feature extractor is improved? To investigate the influence of the backbone architecture on the certified robustness of a deployed model, we change the architecture from ResNet-50 to ViT-Base and plot the target CA curve for ViT-Base in solid blue in Figure 2. We observe that the target ACR obtained by ViT-Base on PACS is higher than the target ACR obtained by ResNet-50 across deformations. **5 A significant improvement of the target certified robustness is achieved by using a stronger backbone architecture.** This result is consistent with the robustness literature (Gowal et al., 2020), where stronger backbones tend to exhibit better robustness, and the domain generalization literature (Gulrajani & Lopez-Paz, 2021), where stronger backbones tend to exhibit better generalization accuracy. *We believe that research on models with better generalization can lead to better certified robustness in unseen domains.*

Q6: Does the generalizability of certified robustness correlate with the perceptual similarity between the source and target domains? In all previous experiments, we considered the average certified accuracy over all possible target domains. We now conduct a more fine-grained study to these target domains individually. We measure the drop in the average certified radius (ΔACR) between the source and target domains with the perceptual similarity between both domains captured by FID (Heusel et al., 2017) and the more robust R-FID (Alfarrar et al., 2022c). To that end, we conduct experiments on PACS where we select one domain as the unseen target and treat the rest as source domains. We train a classifier on the source data and plot the certified accuracy curves against scaling and translation deformations on both the source and target domains in Figure 3 accompanied by ΔACR . We also report the FID and R-FID between the source and target domains. Note that *higher* FID/R-FID indicates *less* similarity of distributions. **6 Perceptual similarity, as captured by FID and R-FID, is not predictive of performance and robustness generalizability.** Surprisingly, the photo domain, which has the highest FID (34.3) and R-FID (87.7) scores, exhibits

the largest certified accuracy generalization. In this case, the ACR for the target domain is higher than the source domain resulting in a negative ΔACR (-0.1 when certifying against translation). The appendix includes experiments with other deformations where we observe similar behavior. *We regard the development of a suitable distribution similarity metric, which better correlates with the level of generalizability, as an important research direction.*

6 REAL-WORLD APPLICATION: MEDICAL IMAGES

To demonstrate the applicability of the DG setup to real-world settings, we investigate the generalization of robustness in medical diagnostics. Data collected by medical imaging techniques, like computed tomography (CT) and magnetic resonance imaging (MRI), is susceptible to noise. This noise includes intensity variations caused by subject movement (Shaw et al., 2019), respiratory motion (Axel et al., 1986), quantum noise associated with X-rays (Hsieh, 1998), and inhomogeneity in the MRI magnetic field (Leemput et al., 1999). Moreover, due to privacy concerns, a model trained in one medical institution should be deployed in another with limited data sharing and retraining (Kaissis et al., 2020; Ziller et al., 2021; Liu et al., 2021). To test the generalization of robustness in this practical setup, we use the DG dataset WILDS CAMELYON17 (Bánda et al., 2019; Koh et al., 2021) to train models on tissue images from four hospitals and evaluate them on images from an unseen hospital. For the first time, in addition to domain generalization, we explore robustness in CAMELYON17. The task in WILDS CAMELYON17 is to predict whether a tissue image contains a cancerous tumor or not.

Adversarial Augmentation for Better Generalizability. We test the generalization accuracy of a standard-trained ($\epsilon = 0$) and a robust ($\epsilon = 8/255$) model by following setup from Section 4. In contrast to the domains studied in Table 2, adversarial training improves generalization to the unseen hospital. While the clean accuracy of the standard model is 94.05%, the clean accuracy of the robust model is 95.28%. This value is even competitive with the target accuracy (95.25%) obtained by the popular DG strategy CORAL (Sun & Saenko, 2016). The robust accuracy also improves from 82.03% to 92.67%. This significant boost in domain generalization can be attributed to the similarity between pixel perturbations and the underlying domain shift in the medical images. *We encourage future works to study different adversarial training methods that go beyond pixel perturbations, and to propose application-specific augmentations for different distribution shifts.*

Certified Robustness. Next, we investigate the generalizability of certified robustness to the unseen hospital. We follow the experimental setup in Section 5 and measure the certified accuracy on the source and target domains. We observe from Figure 4 that some of the certified robustness generalizes to the unseen hospital when evaluated with pixel perturbations and scaling deformations. We include the results for other deformations in the appendix. We note that the drop in certified accuracy to the unseen hospital (given pixel perturbations) is 4 times what we saw in the PACS dataset in Section 5. This is concerning, as many sources of noise affect medical imaging data, so robust medical diagnostics is important for real-world adoption of AI for Health. *We encourage future research to develop better methods to close the target-source gap in certified robustness.*

7 CONCLUSION

We conducted a large scale empirical analysis to study the interplay between adversarial robustness and domain generalization. We deployed adversarial training and randomized smoothing as empirical and certified defenses. We found that both empirical and certified robustness generalize to unseen domains. We further included experiments on a real-world application, where adversarial training benefits both clean and robust accuracy in an unseen domain. Based on our findings, we encourage more research to understand: (i) under which conditions robust training improves the generalization accuracy, and (ii) what methods can improve certified accuracy in unseen domains.

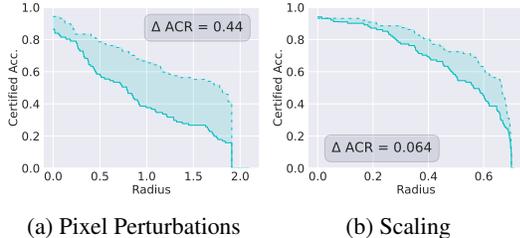


Figure 4: **Certified Robustness in the Medical Domain.** The generalization of robustness against pixel variations in medical images is critical. Yet, there is a significant gap in certified robustness when the DNN is deployed in an unseen hospital.

REFERENCES

- M Alfarra, A Bibi, N Khan, P Torr, and B Ghanem. Deformrs: Certifying input deformations with randomized smoothing. In *Proc. of AAAI Conference on Artificial Intelligence*, 2022a.
- Motasesm Alfarra, Adel Bibi, Philip H. S. Torr, and Bernard Ghanem. Data dependent randomized smoothing. In James Cussens and Kun Zhang (eds.), *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pp. 64–74. PMLR, 01–05 Aug 2022b.
- Motasesm Alfarra, Juan C. Pérez, Anna Frühstück, Philip HS Torr, Peter Wonka, and Bernard Ghanem. On the robustness of quality measures for gans. *arXiv preprint arXiv:2201.13019*, 2022c.
- Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *ICML*, pp. 274–283, 2018.
- L. Axel, R. M. Summers, H. Y. Kressel, and C. Charles. Respiratory effects in two-dimensional fourier transform mr imaging. *Radiology*, 160, 1986. ISSN 00338419. doi: 10.1148/radiology.160.3.3737920.
- Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification tasks to a new unlabeled sample. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- Péter Bándi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong, Quanzheng Li, Farhad Ghazvinian Zanjani, Svitlana Zinger, Keisuke Fukuta, Daisuke Komura, Vlado Ovtcharov, Shenghua Cheng, Shaoqun Zeng, Jeppe Thagaard, Anders B. Dahl, Huangjing Lin, Hao Chen, Ludwig Jacobsson, Martin Hedlund, Melih Çetin, Eren Halici, Hunter Jackson, Richard Chen, Fabian Both, Jörg Franke, Heidi Kusters-Vandeveld, Willem Vreuls, Peter Bult, Bram Van Ginneken, Jeroen Van Der Laak, and Geert Litjens. From detection of individual metastases to classification of lymph node status at the patient level: The camelyon17 challenge. *IEEE Transactions on Medical Imaging*, 38, 2019. ISSN 1558254X. doi: 10.1109/TMI.2018.2867350.
- Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, AISec ’17, pp. 3–14, 2017.
- Fabio M. Carlucci, Antonio D’Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. Domain generalization by solving jigsaw puzzles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and Sungrae Park. Swad: Domain generalization by seeking flat minima. *Advances in Neural Information Processing Systems*, 34, 2021.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pp. 1310–1320. PMLR, 2019.
- Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- Zhun Deng, Linjun Zhang, Kailas Vodrahalli, Kenji Kawaguchi, and James Zou. Adversarial training helps transfer learning via better representations. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.

- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring the landscape of spatial robustness. In *ICML*, 2019.
- Samuel G. Finlayson, John D. Bowers, Joichi Ito, Jonathan L. Zittrain, Andrew L. Beam, and Isaac S. Kohane. Adversarial attacks on medical machine learning. *Science*, 363(6433):1287–1289, 2019.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Robert Geirhos, Kantharaju Narayanappa, Benjamin Mitzkus, Tizian Thieringer, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Partial success in closing the gap between human and machine vision. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.
- Rui Gong, Wen Li, Yuhua Chen, and Luc Van Gool. Dlow: Domain flow for adaptation and generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *CoRR*, abs/1412.6572, 2015.
- Sven Gowal, Chongli Qin, Jonathan Uesato, Timothy Mann, and Pushmeet Kohli. Uncovering the limits of adversarial training against norm-bounded adversarial examples, 2020.
- Shixiang Shane Gu and Luca Rigazio. Towards deep neural network architectures robust to adversarial examples. *CoRR*, abs/1412.5068, 2015.
- Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International Conference on Learning Representations*, 2021.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, 2017.
- Jiang Hsieh. Adaptive streak artifact reduction in computed tomography resulting from excessive x-ray photon noise. *Medical Physics*, 25, 1998. ISSN 00942405. doi: 10.1118/1.598410.
- Georgios A. Kaissis, Marcus R. Makowski, Daniel Rückert, and Rickmer F. Braren. Secure, privacy-preserving and federated machine learning in medical imaging. *Nature Machine Intelligence*, 2, 2020. ISSN 25225839. doi: 10.1038/s42256-020-0186-1.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiko Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton Earnshaw, Imran Haque, Sara M Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. Wilds: A benchmark of in-the-wild distribution shifts. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 5637–5664. PMLR, 18–24 Jul 2021.
- Mathias Lécuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel J. Hsu, and Suman Sekhar Jana. Certified robustness to adversarial examples with differential privacy. *2019 IEEE Symposium on Security and Privacy (SP)*, pp. 656–672, 2019.

- Sungyoon Lee, Woojin Lee, Jinseong Park, and Jaewook Lee. Towards better understanding of training certifiably robust models against adversarial examples. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.
- Koen Van Leemput, Frederik Maes, Dirk Vandermeulen, and Paul Suetens. Automated model-based tissue classification of mr images of the brain. *IEEE Transactions on Medical Imaging*, 18, 1999. ISSN 02780062. doi: 10.1109/42.811270.
- Da Li, Yongxin Yang, Yi Zhe Song, and Timothy M. Hospedales. Deeper, broader and artier domain generalization. volume 2017-October, 2017. doi: 10.1109/ICCV.2017.591.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Learning to generalize: Meta-learning for domain generalization. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Quande Liu, Cheng Chen, Jing Qin, Qi Dou, and Pheng-Ann Heng. Feddg: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1013–1023, June 2021.
- Wang Lu, Jindong Wang, Haoliang Li, Yiqiang Chen, and Xing Xie. Domain-invariant feature exploration for domain generalization. *Transactions on Machine Learning Research*, 2022.
- Xingjun Ma, Yuhao Niu, Lin Gu, Yisen Wang, Yitian Zhao, James Bailey, and Feng Lu. Understanding adversarial attacks on deep learning based medical image analysis systems. *Pattern Recognition*, 110:107332, 2021.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- Jeet Mohapatra, Tsui-Wei Weng, Pin-Yu Chen, Sijia Liu, and Luca Daniel. Towards verifying robustness of neural networks against a family of semantic perturbations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 244–252, 2020.
- Saeid Motiian, Marco Piccirilli, Donald A. Adjeroh, and Gianfranco Doretto. Unified deep supervised domain adaptation and generalization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- A. Tuan Nguyen, Toan Tran, Yarin Gal, and Atilim Gunes Baydin. Domain invariant representation learning with domain density transformations. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.
- Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 372–387, 2016.
- Aditi Raghunathan, Jacob Steinhardt, and Percy Liang. Certified defenses against adversarial examples. In *International Conference on Learning Representations*, 2018.
- Gabriel Pérez S., Juan C. Pérez, Motasem Alfarra, Silvio Giancola, and Bernard Ghanem. 3deformers: Certifying spatial deformations on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15169–15179, June 2022.
- Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 3533–3545. Curran Associates, Inc., 2020.
- Ali Shafahi, Parsa Saadatpanah, Chen Zhu, Amin Ghiasi, Christoph Studer, David Jacobs, and Tom Goldstein. Adversarially robust transfer learning. In *International Conference on Learning Representations*, 2020.

- Richard Shaw, Carole Sudre, Sebastien Ourselin, and M. Jorge Cardoso. Mri k-space motion artefact augmentation: Model robustness and task-specific uncertainty. In M. Jorge Cardoso, Aasa Feragen, Ben Glocker, Ender Konukoglu, Ipek Oguz, Gozde Unal, and Tom Vercauteren (eds.), *Proceedings of The 2nd International Conference on Medical Imaging with Deep Learning*, volume 102 of *Proceedings of Machine Learning Research*, pp. 427–436. PMLR, 08–10 Jul 2019.
- Zheyang Shen, Jiashuo Liu, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards out-of-distribution generalization: A survey, 2021.
- Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation, 2016.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, D. Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *CoRR*, abs/1312.6199, 2014.
- Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *International Conference on Learning Representations*, 2019.
- Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance, 2014.
- Francisco Utrera, Evan Kravitz, N. Benjamin Erichson, Rajiv Khanna, and Michael W. Mahoney. Adversarially-trained deep nets transfer better: Illustration on image classification. In *International Conference on Learning Representations*, 2021.
- Jindong Wang, Culing Lan, Chang Liu, Yidong Ouyang, and Tao Qin. Generalizing to unseen domains: A survey on domain generalization. In Zhi-Hua Zhou (ed.), *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pp. 4627–4635. International Joint Conferences on Artificial Intelligence Organization, 8 2021. Survey Track.
- Mei Wang and Weihong Deng. Deep visual domain adaptation: A survey. *Neurocomputing*, 312: 135–153, 2018.
- Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving adversarial robustness requires revisiting misclassified examples. In *International Conference on Learning Representations*, 2020.
- Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 2958–2969. Curran Associates, Inc., 2020.
- Runtian Zhai, Chen Dan, Di He, Huan Zhang, Boqing Gong, Pradeep Ravikumar, Cho-Jui Hsieh, and Liwei Wang. Macer: Attack-free and scalable robust training via maximizing certified radius. In *International Conference on Learning Representations (ICLR)*, 2020.
- Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. Theoretically principled trade-off between robustness and accuracy. volume 97, pp. 7472–7482. PMLR, 2 2019.
- Xingxuan Zhang, Peng Cui, Renzhe Xu, Linjun Zhou, Yue He, and Zheyang Shen. Deep stable learning for out-of-distribution generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5372–5382, June 2021.
- Kaiyang Zhou, Yongxin Yang, Timothy Hospedales, and Tao Xiang. Learning to generate novel domains for domain generalization. In *European conference on computer vision*, pp. 561–578. Springer, 2020.
- Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. Domain generalization with mixstyle. In *International Conference on Learning Representations*, 2021.
- Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–20, 2022.

Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1): 43–76, 2020.

Alexander Ziller, Dmitrii Usynin, Rickmer Braren, Marcus Makowski, Daniel Rueckert, and Georgios Kaissis. Medical imaging deep learning with differential privacy. *Scientific Reports*, 11, 2021. ISSN 20452322. doi: 10.1038/s41598-021-93030-0.

A EMPIRICAL ROBUSTNESS AND DOMAIN GENERALIZATION

A.1 THE EFFECT OF λ ON ROBUSTNESS GENERALIZATION

How does the parameter λ in equation 5 affect the target robustness? We study the effect of varying λ , which controls the robustness-accuracy trade-off in equation 5. Intuitively, the closer λ is to zero, the more robust the model is. However, this added robustness comes at the cost of clean data accuracy. Similar to 1, which shows the robust accuracies on various datasets with given $\lambda = 0.5$, in Figure 5 we visualize the evaluation results for PACS dataset for $\lambda = 0.1$ and $\lambda = 0.9$. The extreme case of robust only training ($\lambda = 0$) is visualized in Figure 6

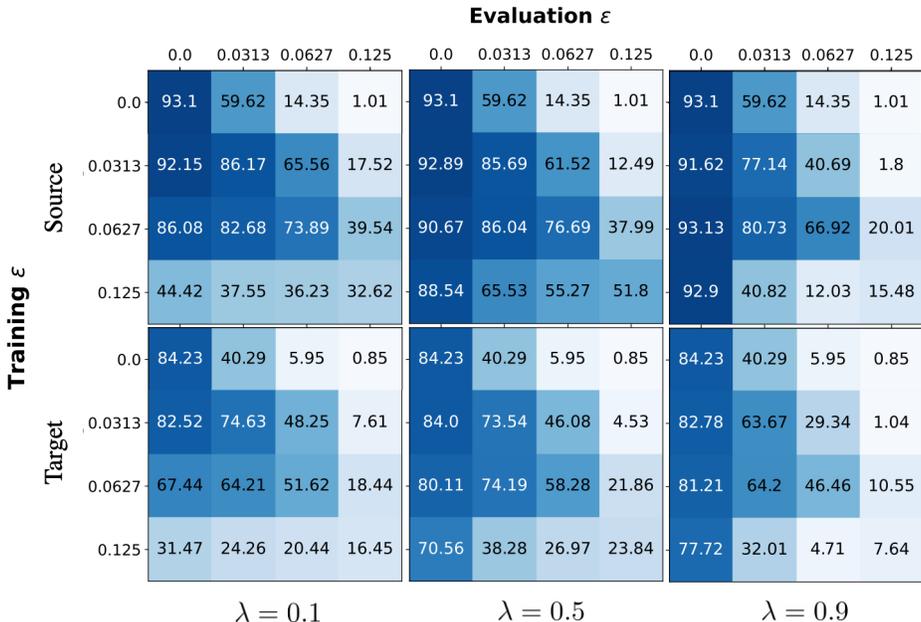


Figure 5: **The Effect of λ on Robustness Generalizability.** As we decrease the value of λ , the network sacrifices the clean accuracy to improve the robust accuracy.

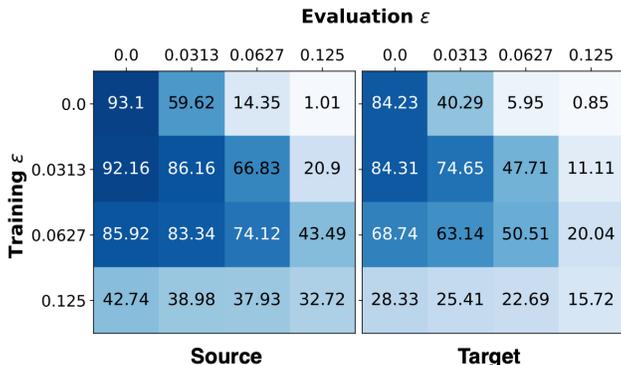


Figure 6: **Robust Training Only ($\lambda = 0$).** In the extreme case, $\lambda = 0$, the network is only trained on adversarial samples from the source domain without any clean source samples. This leads to a sharp drop in the network’s accuracy on both the source and target domains.

A.2 THE GENERALIZATION OF ℓ_2 ROBUSTNESS

Do the paper conclusions about the generalization of robustness hold if we consider ℓ_2 adversarial attacks? We repeat the experiments in Section 4.2 using ℓ_2 adversarial augmentations.

We set $\epsilon = \{0, 0.5, 1.0, 5.0\}$. Each model is trained on one ϵ but is evaluated on all ϵ values. We see from Figure 7 that the paper conclusions also hold when considering ℓ_2 robustness. Following Section 4.2, we answer the following questions:

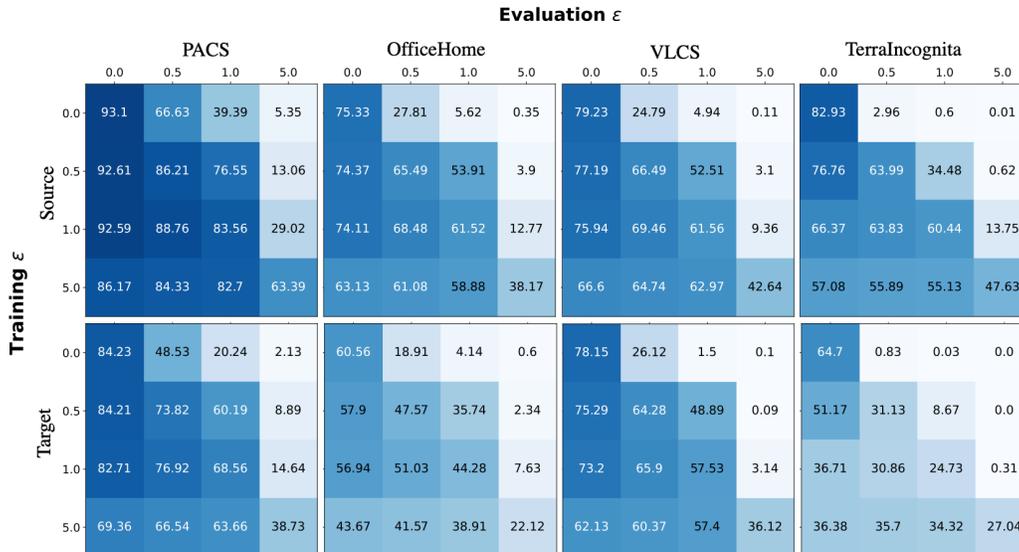


Figure 7: **Evaluation of ℓ_2 Robustness.** The robustness of the source domain (row 1) and the target domain (row 2) follow a similar trend for each dataset. Robustness transfers from the source domains to the target domains, a higher robustness in the source domain is associated with higher robustness in the target domain and vice versa.

Q1: Do adversarially robust models generalize better than their standard-trained counterparts? Again the answer is no. Adversarially trained models tend to experience a drop in generalizability when compared to their standard-trained counterparts.

Q2: Does a higher source-domain robustness correspond to a higher target-domain robustness? As expected, the answer is still yes. As observed across Table 7, when we have a higher robustness in the source domain we consistently observe a higher robustness in the target domain.

Q3: Does the robustness-accuracy trade-off generalize to unseen domains? Yes, similar to what was observed in ℓ_∞ experiments, the robustness-accuracy trade-off exists in unseen domains. Robustness in the target domain comes at the cost of clean accuracy.

A.3 THE EFFECT OF USING A STRONGER EMPIRICAL DEFENSE

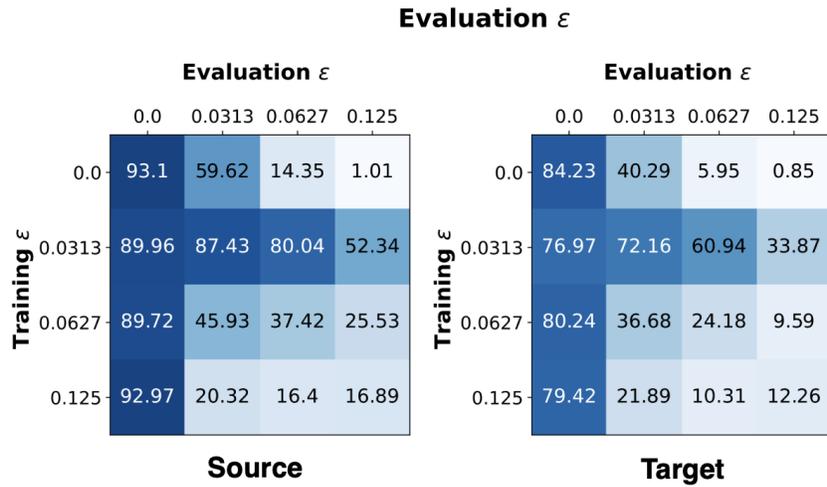


Figure 8: **Robustness Generalizability Using TRADES.** Utilizing a stronger adversarial training method, TRADES, does not provide guarantees towards higher target robustness.

Does the use of a stronger defense, e.g. TRADES (Zhang et al., 2019), improve the generalizability of DNN robustness? When a stronger adversarial robustness method is deployed, in general and unexpectedly we obtain lower robustness generalizability compared to standard adversarial training. We also happen to experience a sharper drop in target accuracy for $\epsilon = 8/255, 16/255$.

B CERTIFIED ROBUSTNESS AND DOMAIN GENERALIZATION

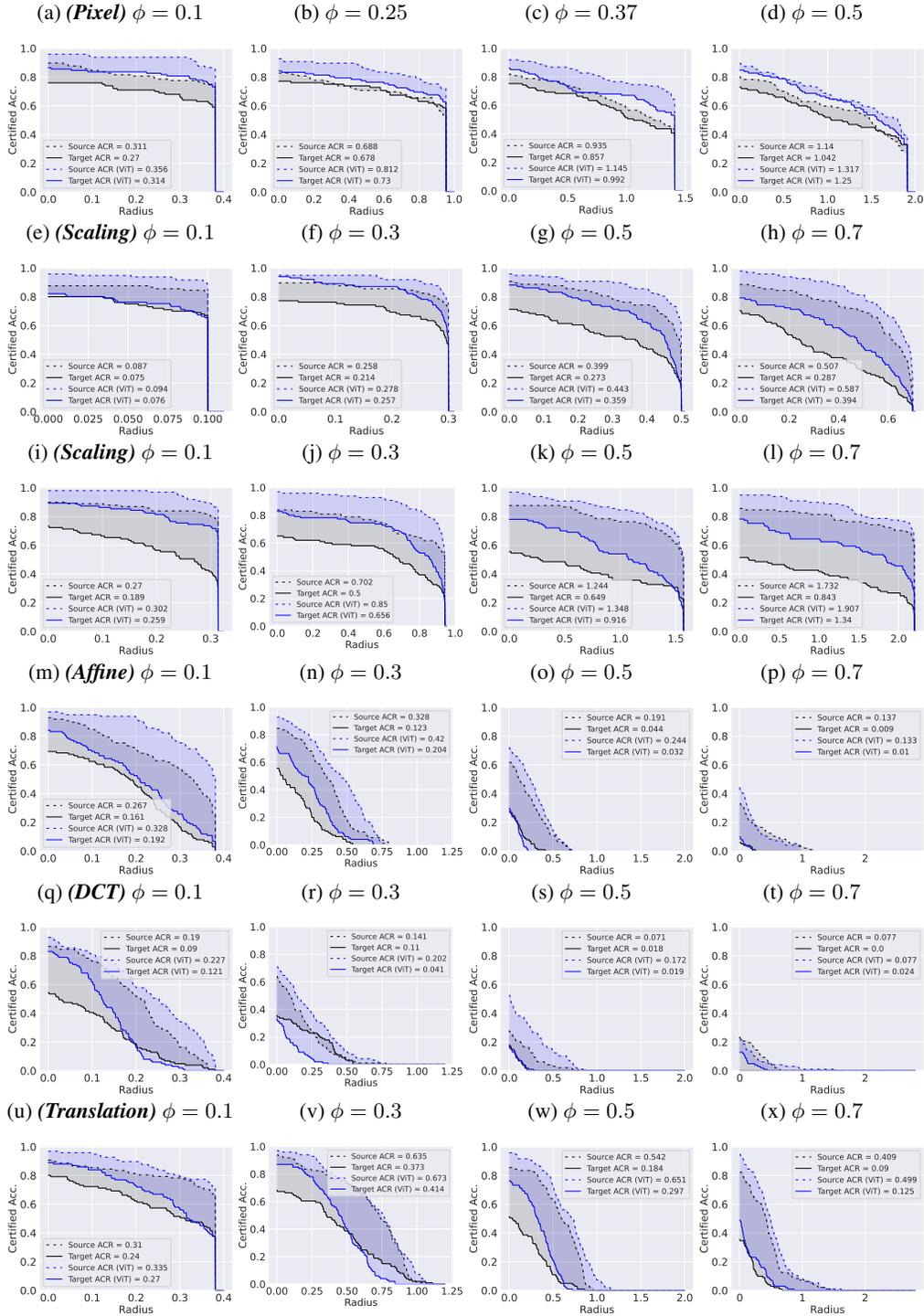


Figure 9: **Effect of Varying the Deformation Parameter ϕ .** We observe that 1) for Pixel Perturbations, Scaling, and Rotation, the higher ϕ gets, the larger the Average Certified Radius (ACR) becomes; and 2) for Affine, DCT, and Translation, high ϕ values can result in low ACRs.

B.1 THE EFFECT OF ϕ ON THE GENERALIZATION OF CERTIFIED ROBUSTNESS

To complement Figure 2, we investigate the behavior when the deformation parameter ϕ varies. Following Section 5.2, we certify ResNet-50 and ViT-Base against pixel perturbations and input deformations in the source and target domains of PACS. We break down each envelope curve in Figure 2 into multiple curves, each representing one choice of ϕ in Eq. 8. We label each curve with the corresponding ϕ value in Figure 9. We observe that the effect of ϕ largely depends on the type of perturbation. On the one hand, for Scaling and Pixel Perturbations, a higher ϕ values corresponds to a larger Average Certified Radius (ACR); on the other hand, for Affine, DCT, and Translation, a higher ϕ values might correspond to a smaller ACR. This is because for the latter group of deformations, a higher ϕ results in a completely deformed image, which hinders the certification ability of the model, even at a small radius. We visualize images from these deformations in Section D. Note that the trends in Figure 2 still stand. Specifically, (1) for ϕ values where there’s a reasonable certified accuracy in the source, that certified accuracy generalizes to the target. Moreover, (2) A stronger architecture (ViT-Base) generally leads to a better source and target certified accuracy.

B.2 DOES VISUAL SIMILARITY CORRELATE WITH ROBUSTNESS GENERALIZABILITY?

We repeat the experiments in Section 5.2, which aim to evaluate the ability of FID/R-FID to predict the generalization of robustness, on pixel perturbations and the following deformations: rotation, affine, and DCT. We observe from Figure 10 that the FID/R-FID values do not predict the level of generalizability of certified robustness, which matches our paper findings.

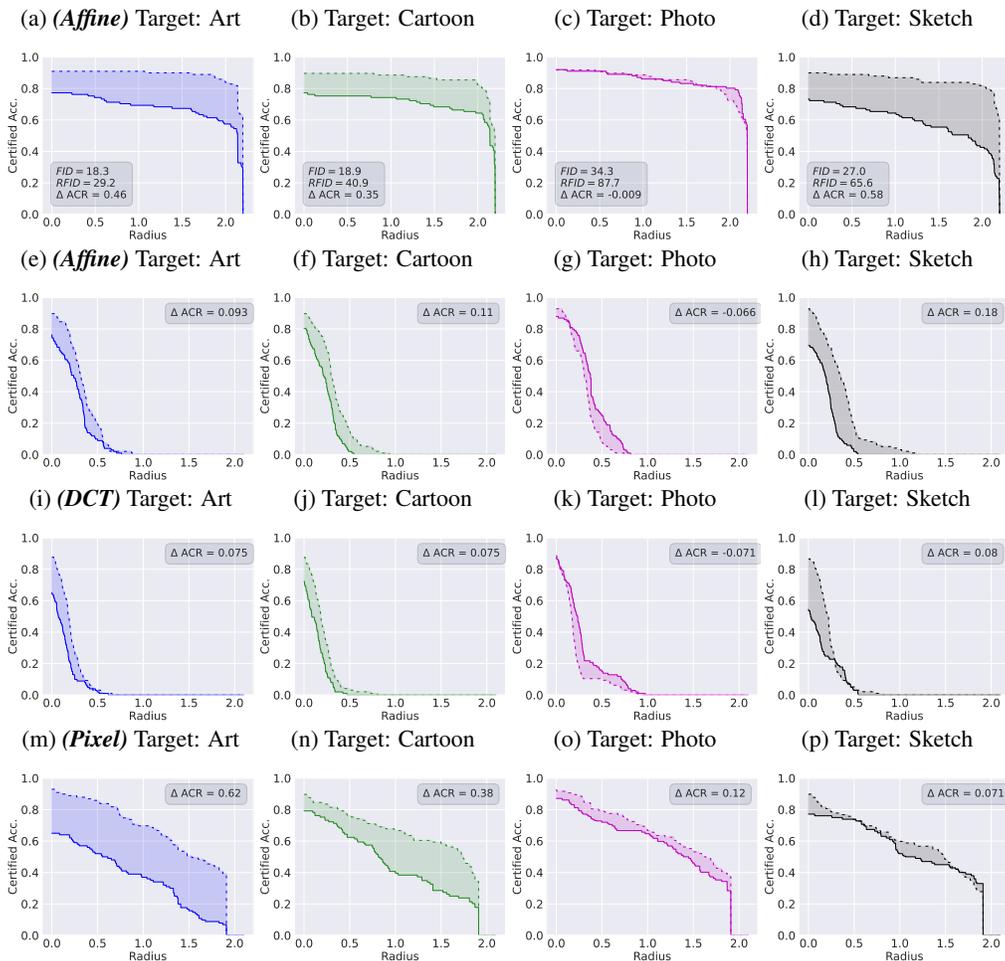


Figure 10: **Does visual similarity correlates with robustness generalizability?** We vary the target distribution and plot the certified accuracy curves for different deformations. The FID/R-FID distances between the source and target distributions are shown in the first row. Visual similarity (FID and R-FID) does not correlate with the level of robustness generalization to the target domain.

C REAL-WORLD APPLICATION: MEDICAL IMAGES

Clean Samples

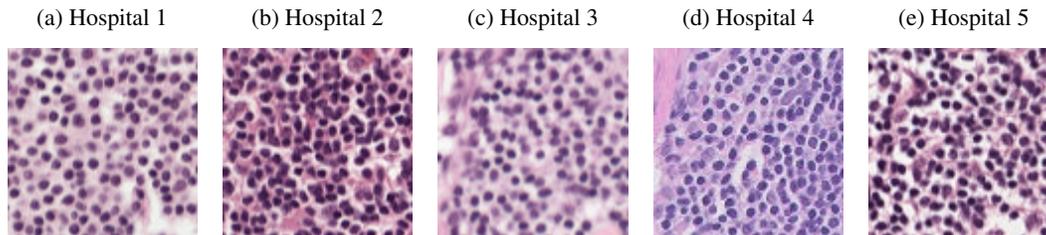


Figure 11: **A visualization of the images taken from the 5 hospitals in Camelyon17.**

We repeat the certified robustness experiments in Section 6 on the following deformations: affine, DCT, translation, and rotation. We observe from Figure 12 that the source-target certification gap is similar for affine, DCT, and translation. However, the certified accuracy curves for rotation are

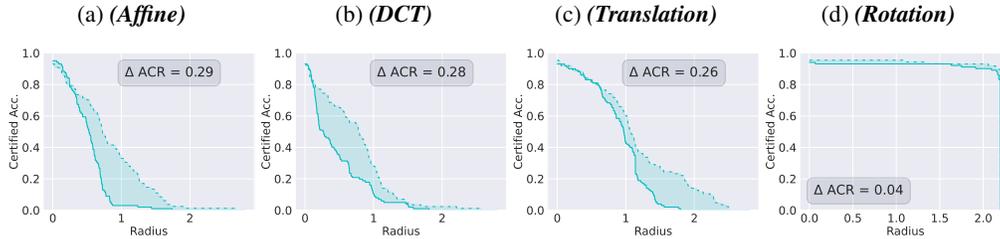


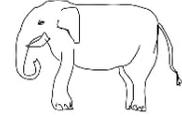
Figure 12: **Does visual similarity correlates with robustness generalizability?** We vary the target distribution and plot the certified accuracy curves for two deformations: scaling and translation. A sample from each distribution is shown in the second row. The FID/R-FID distances between the source distributions and each target are inset in the first row. Visual similarity, measured by FID and R-FID, does not correlate with the level of robustness generalization to the target domain.

different. This makes sense when we consider the way the Camelyon17 dataset is constructed (Bánci et al., 2019; Koh et al., 2021). The dataset includes cropped histopathological images, each of which may contain a tumor tissue in the central 32x32 region. Due to the nature of this construction, rotated versions of the image look similar, which explains why the source and target certified radii remain almost constant. Samples from the Camelyon17 dataset are visualized in Figure 11.

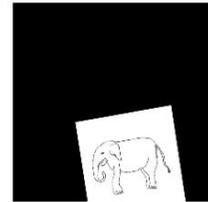
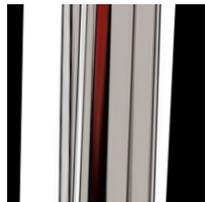
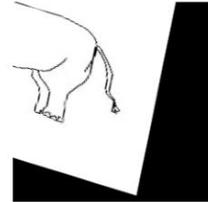
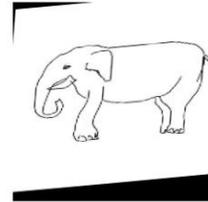
D VISUALIZING THE DOMAIN GENERALIZATION DATASETS

We visualize a few samples from each of the domain generalization datasets we used in the paper. We note that the datasets are diverse in terms of the nature of domain shifts and real-world applicability. Along with the clean samples, we visualize deformed versions of the samples under various values of σ for all the studied deformations.

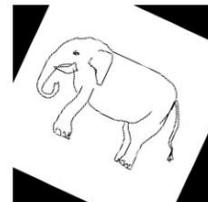
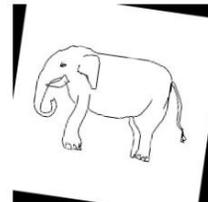
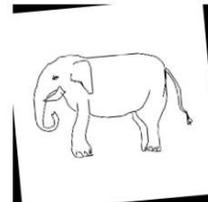
Clean Samples



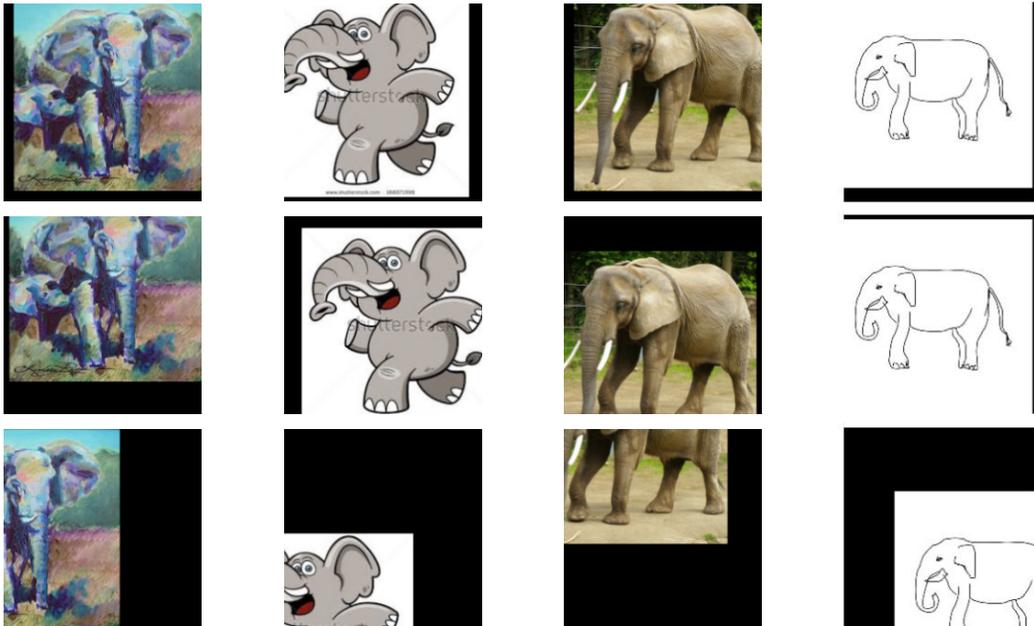
Affine, $\sigma = 0.1, 0.3, 0.5$



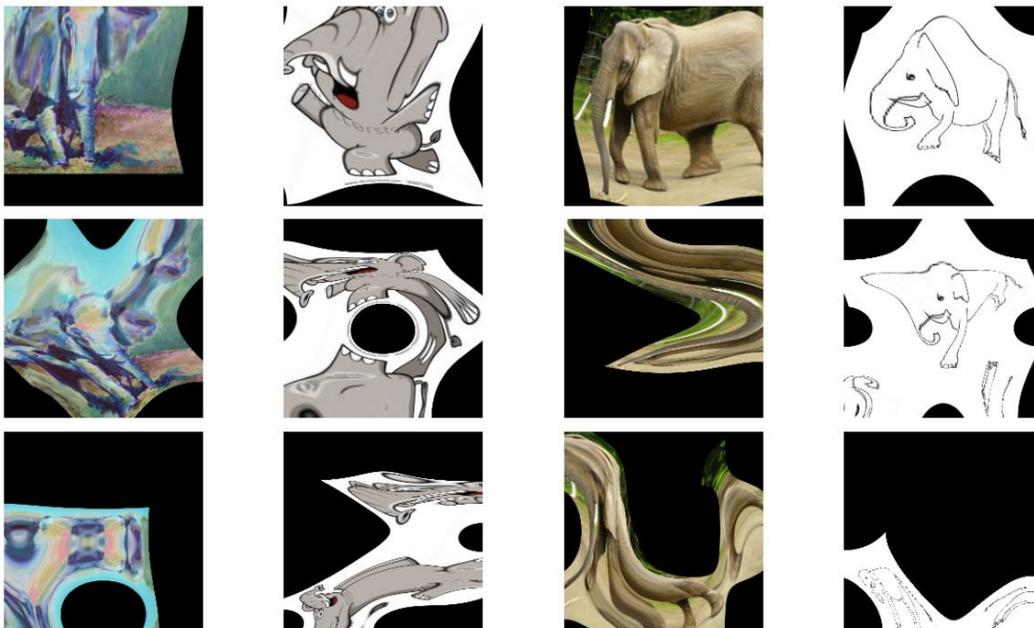
Rotation, $\sigma = 0.1, 0.3, 0.5$



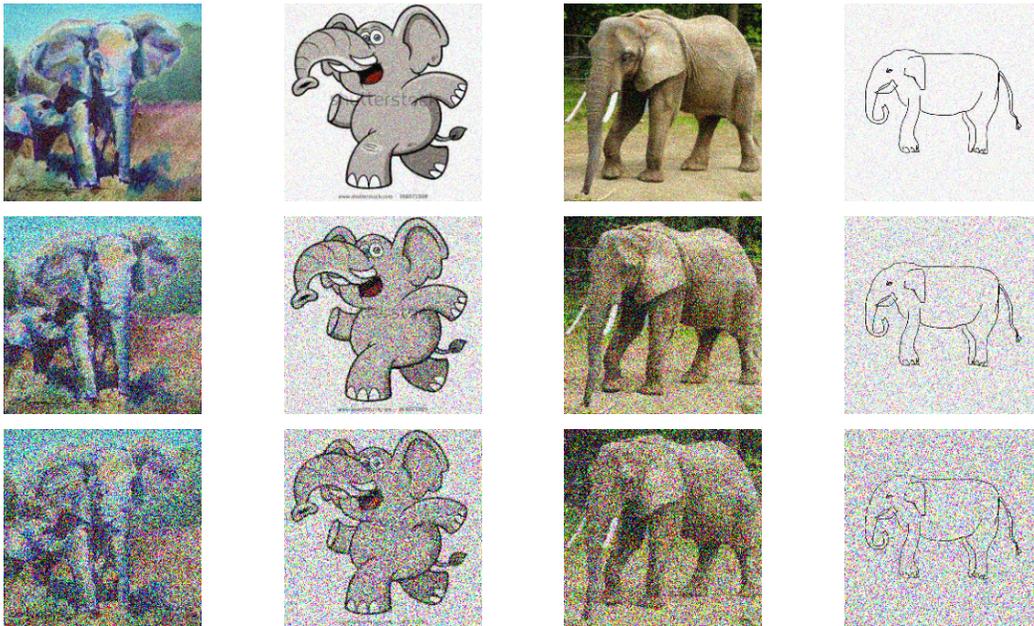
Translation, $\sigma = 0.1, 0.3, 0.5$



DCT, $\sigma = 0.1, 0.3, 0.5$



Pixel Perturbation, $\sigma = 0.1, 0.3, 0.5$



Scaling, $\lambda = 0.1, 0.3, 0.5$

