
Evaluating Worst Case Adversarial Weather Perturbations Robustness

Yihan Wang
UCLA
wangyihan617@gmail.com

Yunhao Ba
UCLA
yhba@ucla.edu

Howard Zhang
UCLA
hwdz15508@g.ucla.edu

Huan Zhang
CMU
huan@huan-zhang.com

Achuta Kadambi
UCLA
achuta@ee.ucla.edu

Stefano Soatto
UCLA
soatto@ucla.edu

Alex Wong
Yale
alex.wong@yale.edu

Cho-Jui Hsieh
UCLA
chohsieh@cs.ucla.edu

Abstract

Several algorithms are proposed to improve the robustness of deep neural networks against adversarial perturbations beyond ℓ_p cases, i.e. weather perturbations. However, evaluations of existing robust training algorithms are over-optimistic. This is in part due to the lack of a standardized evaluation protocol across various robust training algorithms, leading to ad-hoc methods that test robustness on either random perturbations or the adversarial samples from generative models that are used for robust training, which is either uninformative of the worst case, or is heavily biased. In this paper, we identify such evaluation bias in these existing works and propose the first standardized and fair evaluation that compares various robust training algorithms by using physics simulators for common adverse weather effects i.e. rain and snow. With this framework, we evaluated several existing robust training algorithms on two streetview classification datasets (BIC_GSV, Places365) and show the evaluation bias in experiments.

1 Introduction

Adversarial robustness of machine learning models has become an important topic in recent years. For safety-critical applications such as autonomous driving and healthcare systems, it is important to ensure the robustness of models before deploying them into the real world. Most of previous works on adversarial robustness focus on the simplified ℓ_p norm threat model where the worst perturbation is defined within a small ℓ_p ball (Madry et al., 2018; Goodfellow et al., 2015). This simplified assumption ensures the perturbation is imperceptible and makes it easier to design attack and defense algorithms. However, in practice there are many other semantic-preserved or natural perturbations that are not ℓ_p norm bounded and difficult to be mathematically defined, such as adversarial shadows (Zhong et al., 2022), adversarial rains (Zhai et al., 2020), and physical adversarial T-shirts (Xu et al., 2020). Since collecting real-world adversarial examples is infeasible, a natural solution to improve the natural adversarial robustness is through learning a perturbation set that introduces weather effects such as rain or snow into the original clean images, and then conduct adversarial training with the ℓ_p ball constraint replaced by the general perturbation set (Wong & Kolter, 2020; Robey et al., 2020). However, the lack of natural adversarial examples also makes evaluating the robustness of such a perturbation set with real images impossible. To bypass this difficulty, Wong & Kolter

(2020) reuses the pretrained generative model to produce adversarial examples, which is also used in adversarial training. In addition, Robey et al. (2020) deploys out-of-distribution (OOD) perturbations in evaluation, which limits the exploration in the perturbation space. The ad-hoc nature of these evaluations makes it difficult to compare existing robust training algorithms. This stems from a lack of a standardized and fair evaluation protocol for natural robust training algorithms, where current evaluation methods potentially lead to unreliable robustness metrics.

To address the above issue, we propose an evaluation framework for robust training algorithms based on physics-based weather simulators. The framework is designed in a differentiable manner to enable the generation of worst-case adversarial examples during evaluation. The adversarial examples generated by physical simulators can give a fair evaluation by using independent random simulators both in training and evaluation. In experiments, we show that several existing models are robust to the perturbations from their own generative models as expected, but are not robust to our adversarial attacks based on physical simulators.

Our main contributions are summarized as follows:

- We propose the first standard protocol for evaluating robust training algorithms for machine learning models against weather perturbations. We do so by leveraging physics-based simulators to model adverse weather effects – yielding fair comparisons between algorithms for learning perturbation sets beyond ℓ_p norm and robust training.
- We demonstrate that existing robust training methods are over-optimistic in their claims due to their ad-hoc evaluations, which are either uninformative of worst case perturbations or favorable towards the particular training algorithm.

2 Background

We consider the robust training problem where we are given a clean dataset and a perturbed dataset, which is normally the case when we want to train a model robust against natural perturbations where the perturbation set is difficult to be defined mathematically. For example, we can collect a dataset of sunny images $(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n)$ and a dataset of rainy i.e. perturbed images $(\mathbf{x}'_0, \mathbf{x}'_1, \dots, \mathbf{x}'_m)$. Here the clean dataset is sampled from the clean distribution $\mathbf{x} \in \mathbb{R}^d \sim p(\mathbf{x})$. Each perturbed example \mathbf{x}' in the perturbed dataset is corresponding to a clean image \mathbf{x} and $\mathbf{x}' = g(\mathbf{x}, \mathbf{z})$ where \mathbf{z} is the feature vector for this perturbation and g is the perturbation function. In this work, we focus on rain and snow perturbations where we can write the perturbation function as $\mathbf{x}' = g(\mathbf{x}, \mathbf{z}) = \mathbf{x} + \delta = \mathbf{x} + g'(\mathbf{x}, \mathbf{z})$. Here g' is the function that generates the perturbation mask δ from the perturbation feature \mathbf{z} . Here \mathbf{z} can be limited within a pre-defined perturbation set $\mathbf{z} \in \Delta$ which limits the perturbation δ from being too strong to alter the semantic meaning of the clean image \mathbf{x} .

A classifier trained with a (robust) training algorithm $f_\theta(\mathbf{x}) : \mathbb{R}^d \rightarrow C$ parameterized by θ is considered as safe for an example \mathbf{x} if we cannot find an adversarial example with an adversarial attack algorithm. In adversarial attack, we find the worst-case examples:

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{z} \in \Delta} l(f_\theta(\mathbf{x}'), y). \tag{1}$$

Here $l(\cdot, y)$ is the standard cross entropy loss where y is the ground-truth label.

If $f_\theta(\hat{\mathbf{x}}) = y$, f_θ is safe for \mathbf{x} under adversarial attack. In this work, we evaluate the robustness of a given model f_θ by evaluating the percentage of safe examples in the test set.

3 Evaluation Framework for Robust Training Algorithms

In this section, we describe our evaluation framework for robust training algorithms. We illustrate the overview of this framework in Figure 1. Given a robust training algorithm, we run the algorithm with the same datasets generated by our framework and then evaluate the result robust model with the same environment.

Differentiable physics-based weather simulator: The key module of this framework is the differentiable simulation engine which enables us to apply adversarial attacks directly to evaluate the

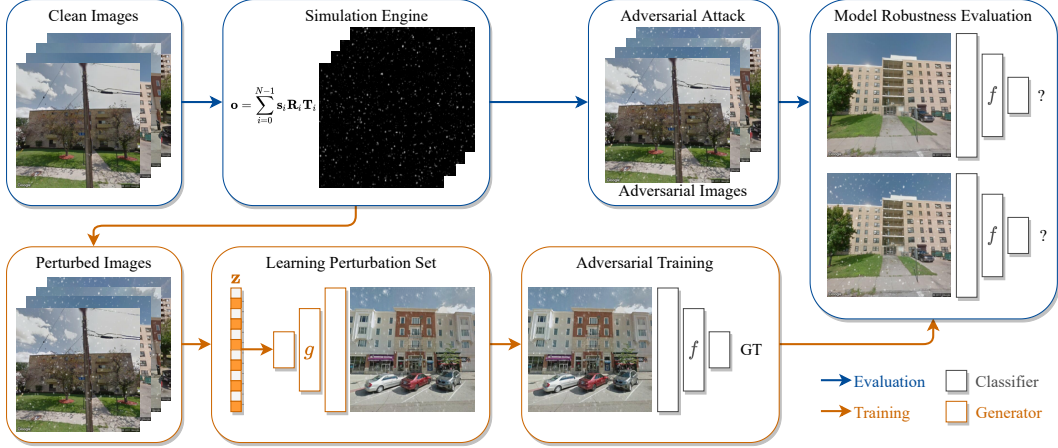


Figure 1: An overview of our proposed evaluation framework.

robustness of downstream models without the need for ad-hoc generative models used by previous works (Wong & Kolter, 2020; Robey et al., 2020).

The simulator generates weather perturbations in a two-stage process: (1) particles rendering, and (2) particles aggregation.

Falling particles like rain drops and snowflakes are rendered according to physical parameters such as particle size, scene illumination, and camera settings. For rain drops, we follow Garg & Nayar (2006). For snowflakes, we generate each snowflake from a Gaussian noise image and control the size, shape, and direction with threshold, standard deviation, and a motion blur, respectively. These diverse rain drops and snowflakes are pre-generated with random parameters and stored in a database.

At the second stage, we randomly sample N particles from the particle database and aggregate them to form a weather mask. To make it differentiable, each particle image denoted by \mathbf{s}_i is parameterized by a translation matrix \mathbf{T}_i and a rotation matrix \mathbf{R}_i . We can generate a weather effect mask with N particles via:

$$\mathbf{o} = \sum_{i=0}^{N-1} \mathbf{s}_i \mathbf{R}_i \mathbf{T}_i. \quad (2)$$

Then the mask can be merged with the clean image to produce a perturbed image $\mathbf{x}' = \mathbf{o} + \mathbf{x}(1 - \mathbf{o})$ where \mathbf{x} is the clean image. This differentiable aggregation enables us to generate adversarial examples directly in the simulator:

$$\mathbf{R}_i^*, \mathbf{T}_i^* = \arg \max_{\mathbf{R}_i, \mathbf{T}_i} l(f(\mathbf{x}'), y),$$

where \mathbf{R}_i and \mathbf{T}_i are constrained in some given perturbation set. The constrained maximization problem can be solved with PGD attack (Madry et al., 2018). We call the adversarial attack with our differentiable simulator as simulation attack.

Evaluation framework: To evaluate robust training algorithms, we have a clean dataset which is the same for all training algorithms. Then we use the simulator described earlier to generate a random perturbed dataset with random transformation and rotation matrix. The clean and perturbed datasets are then fed into the training algorithm, where potentially the training algorithm includes a module to learn perturbation set and then conduct adversarial training within this learned perturbation set. After obtaining the trained robust model, we use the differentiable simulator to evaluate the adversarial natural robustness of this model with adversarial attack on the simulation parameters \mathbf{R}_i and \mathbf{T}_i .

4 Experiments

In this section, we evaluate some existing robust training algorithms with our proposed evaluation framework and compare the results with their own evaluations. We will show that evaluations in existing work have bias which prefers the models trained with its own generative models.

Table 1: Ranks and robust error (R.E.) of different training methods under different evaluations for BIC_GSV (Kang et al., 2018) with rainy effects. Method with the lowest robust error is ranked as 1. “Random” is evaluated with randomly generated perturbed images, VRGNet (Wang et al., 2021), CVAE (Wong & Kolter, 2020) and MUNIT (Robey et al., 2020) are evaluated with adversarial attack based on different generators and “Attack” is evaluated with simulation attack.

Training methods	Random		VRGNet		CVAE		MUNIT		Attack	
	Rank	R.E.	Rank	R.E.	Rank	R.E.	Rank	R.E.	Rank	R.E.
Clean Training	4	0.6424	5	0.7634	5	0.4893	5	0.8027	5	0.6890
Augmented Training	3	0.4815	3	0.5656	2	0.4592	3	0.5928	3	0.5136
VRGNet AT (Wang et al., 2021)	1	0.4587	1	0.4767	3	0.4742	1	0.5063	1	0.4665
CVAE AT (Wong & Kolter, 2020)	5	0.6113	4	0.7498	3	0.4742	4	0.7430	4	0.6477
MUNIT AT (Robey et al., 2020)	2	0.4728	2	0.5228	1	0.4397	1	0.5063	2	0.4893

4.1 Settings

Datasets: In our experiments, we choose two streetview classification datasets which are suitable for weather effects. The BIC_GSV (Kang et al., 2018) dataset contains streetview figures extracted from Google StreetView labeled as 8 classes. Another dataset is the commonly used Places365 dataset (Zhou et al., 2017) for scene recognition which includes 365 scene categories. We picked a subset with 7 outdoor scenes from Places365 as the dataset used in our experiments. We also picked another 4 classes for our unpaired dataset. A list of categories chosen in the subset can be found in the Appendix. For training the unpaired generator models, we generate the perturbed dataset from 4 additional classes for Places dataset. In all of our experiments, the input images are resized to 256×256 pixels. The number of particles in simulation is set to be 2,000.

Robust training algorithm baselines: We include three robust training algorithms in our evaluation experiments. VRGNet (Wang et al., 2021) proposes a Bayesian weather generation model to generate rain/snow masks with an input latent code which can be used in adversarial training by attacking the latent code. Robey et al. (2020) implements MUNIT (Huang et al., 2018) as the generator model which maps clean images to naturally perturbed images with a latent code which can also be used in adversarial training. For experiments related to MUNIT, we implemented a modified version of MUNIT for weather perturbations. The architecture of MUNIT used in our experiments can be found in the Appendix. Wong & Kolter (2020) learns the perturbation set with a CVAE model conditioned on a latent code. Both VRGNet and CVAE require paired datasets and MUNIT supports unpaired datasets. The three generators can then be used in adversarial training to train robust downstream models. Details of training settings can be found in the Appendix.

4.2 Evaluation Results

We evaluate and compare the robust error of five models trained with different methods under the five evaluations for BIC_GSV (Kang et al., 2018) under rain perturbations. Clean training is trained with only the clean dataset. Augmented Training is augmented with randomly perturbed images from the simulator. For the five different evaluation methods, we use 10-step PGD attack in all methods that requires adversarial attack. In Table 1, we list the ranks of different training methods under different evaluations where the model with the lowest robust error is ranked as 1. We also report the exact numbers in evaluation.

As illustrated in Table 1, if we train and test a model using the same generator, the model can overfit the perturbation set defined by the generator during training. For example, when evaluating the model with adversarial attacks using VRGNet or MUNIT, the model trained with the corresponding AT has the best adversarial robust accuracy. This result validates our motivation to propose a standard evaluation protocol for robust training algorithms.

5 Conclusion

In the paper, we propose the first standard evaluation protocol for robust training algorithms under perturbation sets beyond ℓ_p . This addresses a gap in the community as evaluations were previously performed using ad-hoc generative models that were also used by the robust training algorithms, which can lead to unreliable evaluation.

References

- Peter C Barnum, Srinivasa Narasimhan, and Takeo Kanade. Analysis of rain and snow in frequency space. *International journal of computer vision*, 86(2):256–274, 2010.
- Battista Biggio, Iginio Corona, Davide Maiorca, Blaine Nelson, Nedim Šrđić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 387–402. Springer, 2013.
- Dan A Calian, Florian Stimberg, Olivia Wiles, Sylvestre-Alvise Rebuffi, Andras Gyorgy, Timothy Mann, and Sven Gowal. Defending against image corruptions through adversarial augmentations. *arXiv preprint arXiv:2104.01086*, 2021.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *Security and Privacy (SP), 2017 IEEE Symposium on*, pp. 39–57, 2017.
- Wei-Ting Chen, Hao-Yu Fang, Jian-Jiun Ding, Cheng-Che Tsai, and Sy-Yen Kuo. Jstasr: Joint size and transparency-aware snow removal algorithm based on modified partial convolution and veiling effect removal. In *European Conference on Computer Vision*, pp. 754–770. Springer, 2020.
- Wei-Ting Chen, Hao-Yu Fang, Cheng-Lin Hsieh, Cheng-Che Tsai, I Chen, Jian-Jiun Ding, Sy-Yen Kuo, et al. All snow removed: Single image desnowing algorithm using hierarchical dual-tree complex wavelet representation and contradict channel loss. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4196–4205, 2021.
- Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9185–9193, 2018.
- Ranjie Duan, Xingjun Ma, Yisen Wang, James Bailey, A Kai Qin, and Yun Yang. Adversarial camouflage: Hiding physical-world attacks with natural styles. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1000–1008, 2020.
- Ranjie Duan, Xiaofeng Mao, A Kai Qin, Yuefeng Chen, Shaokai Ye, Yuan He, and Yun Yang. Adversarial laser beam: Effective physical-world attack to dnns in a blink. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16062–16071, 2021.
- Kshitiz Garg and Shree K Nayar. Photorealistic rendering of rain streaks. *ACM Transactions on Graphics (TOG)*, 25(3):996–1002, 2006.
- Kshitiz Garg and Shree K. Nayar. Vision and rain. *International Journal of Computer Vision*, 75(1): 3–27, 2007.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015. URL <http://arxiv.org/abs/1412.6572>.
- Shirsendu Sukanta Halder, Jean-François Lalonde, and Raoul de Charette. Physics-based rendering for improving robustness to rain. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10203–10212, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016. doi: 10.1109/CVPR.2016.90. URL <https://doi.org/10.1109/CVPR.2016.90>.
- Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781*, 2019.
- Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8340–8349, 2021.

- Xiaowei Hu, Chi-Wing Fu, Lei Zhu, and Pheng-Ann Heng. Depth-attentional features for single-image rain removal. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8022–8031, 2019.
- Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 172–189, 2018.
- Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. Certified robustness to adversarial word substitutions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4129–4142, Hong Kong, China, 2019. doi: 10.18653/v1/D19-1423. URL <https://aclanthology.org/D19-1423>.
- Jian Kang, Marco Körner, Yuanyuan Wang, Hannes Taubenböck, and Xiao Xiang Zhu. Building instance classification using street view images. *ISPRS journal of photogrammetry and remote sensing*, 145:44–59, 2018.
- Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- Juncheng Li, Frank Schmidt, and Zico Kolter. Adversarial camera stickers: A physical camera-based attack on deep learning systems. In *International Conference on Machine Learning*, pp. 3896–3904. PMLR, 2019a.
- Ruoteng Li, Loong-Fah Cheong, and Robby T. Tan. Heavy rain image restoration: Integrating physics model and conditional adversarial learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1633–1642, 2019b.
- Yu Li, Robby T. Tan, Xiaojie Guo, Jiangbo Lu, and Michael S. Brown. Rain streak removal using layer priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2736–2744, 2016.
- Yun-Fu Liu, Da-Wei Jaw, Shih-Chia Huang, and Jenq-Neng Hwang. Desnownet: Context-aware deep network for snow removal. *IEEE Transactions on Image Processing*, 27(6):3064–3073, 2018.
- 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. URL <https://openreview.net/forum?id=rJzIBfZAb>.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2574–2582, 2016. doi: 10.1109/CVPR.2016.282. URL <https://doi.org/10.1109/CVPR.2016.282>.
- Fabio Pizzati, Pietro Cerri, and Raoul de Charette. Model-based occlusion disentanglement for image-to-image translation. In *European conference on computer vision*, pp. 447–463. Springer, 2020a.
- Fabio Pizzati, Raoul de Charette, Michela Zaccaria, and Pietro Cerri. Domain bridge for unpaired image-to-image translation and unsupervised domain adaptation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2990–2998, 2020b.
- Alexander Robey, Hamed Hassani, and George J Pappas. Model-based robust deep learning: Generalizing to natural, out-of-distribution data. *arXiv preprint arXiv:2005.10247*, 2020.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- Zhiqiang Shen, Mingyang Huang, Jianping Shi, Xiangyang Xue, and Thomas S Huang. Towards instance-level image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3683–3692, 2019.

- Aman Sinha, Hongseok Namkoong, Riccardo Volpi, and John Duchi. Certifying some distributional robustness with principled adversarial training. *arXiv preprint arXiv:1710.10571*, 2017.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014. URL <http://arxiv.org/abs/1312.6199>.
- Maxime Tremblay, Shirsendu Sukanta Halder, Raoul De Charette, and Jean-François Lalonde. Rain rendering for evaluating and improving robustness to bad weather. *International Journal of Computer Vision*, 129(2):341–360, 2021.
- Alexander Von Bernuth, Georg Volk, and Oliver Bringmann. Simulating photo-realistic snow and fog on existing images for enhanced cnn training and evaluation. In *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 41–46. IEEE, 2019.
- Hong Wang, Zongsheng Yue, Qi Xie, Qian Zhao, Yefeng Zheng, and Deyu Meng. From rain generation to rain removal. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14791–14801, 2021.
- Yoann Weber, Vincent Jolivet, Guillaume Gilet, and Djamchid Ghazanfarpour. A multiscale model for rain rendering in real-time. *Computers & Graphics*, 50:61–70, 2015.
- Yanyan Wei, Zhao Zhang, Yang Wang, Mingliang Xu, Yi Yang, Shuicheng Yan, and Meng Wang. Deraincyclegan: Rain attentive cyclegan for single image deraining and rainmaking. *IEEE Transactions on Image Processing*, 30:4788–4801, 2021.
- Eric Wong and J Zico Kolter. Learning perturbation sets for robust machine learning. *arXiv preprint arXiv:2007.08450*, 2020.
- Eric Wong, Frank Schmidt, and Zico Kolter. Wasserstein adversarial examples via projected sinkhorn iterations. In *International Conference on Machine Learning*, pp. 6808–6817. PMLR, 2019.
- Chaowei Xiao, Bo Li, Jun-Yan Zhu, Warren He, Mingyan Liu, and Dawn Song. Generating adversarial examples with adversarial networks. *arXiv preprint arXiv:1801.02610*, 2018.
- Kaidi Xu, Gaoyuan Zhang, Sijia Liu, Quanfu Fan, Mengshu Sun, Hongge Chen, Pin-Yu Chen, Yanzhi Wang, and Xue Lin. Adversarial t-shirt! evading person detectors in a physical world. In *European conference on computer vision*, pp. 665–681. Springer, 2020.
- Liming Zhai, Felix Juefei-Xu, Qing Guo, Xiaofei Xie, Lei Ma, Wei Feng, Shengchao Qin, and Yang Liu. It’s raining cats or dogs? adversarial rain attack on dnn perception. *arXiv preprint arXiv:2009.09205*, 2020.
- He Zhang and Vishal M Patel. Density-aware single image de-raining using a multi-stream dense network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 695–704, 2018.
- He Zhang, Vishwanath Sindagi, and Vishal M Patel. Image de-raining using a conditional generative adversarial network. *IEEE transactions on circuits and systems for video technology*, 30(11): 3943–3956, 2019.
- Yiqi Zhong, Xianming Liu, Deming Zhai, Junjun Jiang, and Xiangyang Ji. Shadows can be dangerous: Stealthy and effective physical-world adversarial attack by natural phenomenon. *arXiv preprint arXiv:2203.03818*, 2022.
- Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

A Appendix

A.1 Related Work

Adversarial robustness of beyond ℓ_p robustness: Neural Networks are shown to be vulnerable to perturbations that are imperceptible to humans. Several papers (Biggio et al., 2013; Carlini & Wagner, 2017; Goodfellow et al., 2015) have shown that neural networks can be attacked by small perturbations bounded by ℓ_1, ℓ_2, ℓ_p balls. For instance, Szegedy et al. (2014) shows that such perturbations can alter the output of a classification network. Goodfellow et al. (2015) introduces the fast gradient sign method (FGSM). Dong et al. (2018); Kurakin et al. (2016); Madry et al. (2018) extend FGSM to iterative optimization to boost its performance. Moosavi-Dezfooli et al. (2016) finds the minimal perturbation to alter the predicted class while Madry et al. (2018) proposed projected gradient descent (PGD) to find the worst-case perturbations.

Beyond ℓ_p robustness, some recent papers extend the perturbation sets to settings that can preserve semantic meaning. Some can be well-defined mathematically such as Wasserstein robustness (Wong et al., 2019), distributional shifts (Sinha et al., 2017; Sagawa et al., 2019) and word substitution (Jia et al., 2019) in texts. Others can not be well-defined but perturbed datasets can be generated or collected, like adversarial shadows (Zhong et al., 2022), adversarial rains (Zhai et al., 2020) and some physical adversarial perturbations (Duan et al., 2020, 2021; Li et al., 2019a; Xu et al., 2020).

Robust training algorithms for natural robustness: For general natural robustness settings, attack methods like PGD (Madry et al., 2018) and FGSM (Goodfellow et al., 2015) cannot be directly applied within the perturbation set. Some works use generative adversarial networks (GANs) to learn the perturbation set and then generate adversarial examples during adversarial training (Xiao et al., 2018; Wong & Kolter, 2020; Robey et al., 2020). Another line of work uses random perturbations to improve the out-of-distribution robustness (Hendrycks et al., 2021, 2019; Calian et al., 2021). However, these solutions cannot bypass the lack of natural adversarial examples in evaluation. Wong & Kolter (2020) reuses the learned generators to generate adversarial examples, where the learned generators is not guaranteed to match the perturbation set under the worst-case perturbations, and evaluating and training the models with the same generative model can leads to fake robustness in evaluation. Robey et al. (2020) uses out-of-distribution (OOD) perturbations, which limits the exploration in perturbation space.

Weather simulation: The appearance of falling particles in rain and snow are highly complicated and can be affected by multiple factors, such as the particle properties, camera configurations, and environmental illumination (Garg & Nayar, 2007; Barnum et al., 2010). Copious amounts of research has been conducted to simulate the weather dynamics based on different principles, including raindrop oscillation models (Garg & Nayar, 2006; Li et al., 2016), frequency domain analysis (Barnum et al., 2010; Weber et al., 2015), and depth dependent formulations (Hu et al., 2019; Li et al., 2019b; Halder et al., 2019; Von Bernuth et al., 2019; Tremblay et al., 2021). More recently, several data-driven deep learning techniques have also been developed for weather effect simulation (Shen et al., 2019; Pizzati et al., 2020a,b; Wang et al., 2021; Wei et al., 2021). Other than these, the community also resorts to some existing image editing software (e.g. PhotoShop) for weather simulation (Liu et al., 2018; Zhang & Patel, 2018; Zhang et al., 2019; Chen et al., 2020, 2021).

A.2 Experiment Details

Datasets For Places dataset, we picked 7 categories from the 365 categories in Places365 Zhou et al. (2017) dataset as the training and test datasets for classifier training, including: apartment_building-outdoor, church-outdoor, garage-outdoor, general_store-outdoor, house, library-outdoor and office_building. For generator training which supports unpaired dataset, we picked 4 additional categories to generate the perturbed dataset, which includes art_gallery, entrance_hall, golf_course and yard.

Classifier architectures We use ResNet34 He et al. (2016) as our classifier architecture. In our classifier training, we initialize the Resnet34 classifier with pretrained weights on Imagenet provided by `cnn_finetune`¹ and the final linear layer is randomly initialized.

¹<https://github.com/creafz/pytorch-cnn-finetune>



Figure 2: Examples for random simulation and simulation attack for rainy effects. (a) is the original clean image sampled from BIC_GSV Kang et al. (2018) dataset.; (b) is produced by random rain simulation; (c) is an adversarial example for a ResNet34 classifier pretrained with clean training; (d) is the difference between (b) and (c).

Generator architectures We implement MUNIT based on pytorch-CycleGAN-and-pix2pix.² We denote a convolutional layer with padding 1, $k 4 \times 4$ filters and stride s as Ck-s. Then the discriminators can be represented as C64-2 C128-2 C256-2 C512-2 C512-1 C1-1.

For the generators, we let x -Rk denote x ResNet blocks, each of which contains two 3×3 convolutional blocks with k filters, TCk denote a transpose convolutional block with k filters of size 3, stride 2, padding 1, output padding 1, SCk denote a convolutional block with k filters of size 3, stride 2, padding 1 and Ck denote a convolutional block with k filters of size 7 and padding 0. Then the mask generator can be represented as 9-R256 TC256 TC256 C3. The generator from perturbed image to clean image can be represented as C64 SC128 SC256 9-R256 TC256 TC256 C3. In our MUNIT implementation, we insert the latent code by concatenating into the hidden layer of the generator from clean images to perturbed images. The generator from clean image to perturbed image can be represented as C64 SC128 SC256 9-R257 TC257 TC257 C3 where the latent code which is transformed by a full-connected layer is concatenated with the output of SC256. For our model and MUNIT, we also have a simple convolutional neural network to reconstruct the latent code which takes the output of 9-R256 as the input. The architecture can be represented C16 C32 C64 L128 L128 where Ck is a convolutional layer with $k 4 \times 4$ filters, stride 2 and padding 1 and Lk is a fully-connected layer with output dimension k .

For VRGNet and CVAE, we run the code provided by authors and we scale the architectures defined in the VRGNet code to fit the image size in our dataset where the input is randomly cropped to 224×224 in training. The weather mask generator architecture in VRGNet is L12544 TC128-4-2-1 TC64-4-2-1 TC32-4-2-1 TC3-8-4-2 and the encoder architecture is C32-8-4-2 C64-4-2-1 C128-4-2-1 C256-4-2-1 L256 where Lk is a fully-connected layer with output dimension k and TCk-f-s-p and Ck-f-s-p are a transposed convolutional and convolutional layers with k filters of size f , stride s and padding p .

Training settings for baseline robust training algorithms The number of dimensions for the latent code in VRGNet, MUNIT is set at 128. MUNIT is trained with a learning rate of $2e-4$ for 5 epochs. For VRGNet, we ran the code provided by the authors for 100 epochs on our datasets with default hyper-parameters. For the CVAE model, we ran the code provided by the authors with the same configuration as the Multi-illumination dataset. For VRGNet and CVAE, we scaled the encoder and decoder to fit the image size in our datasets. Details of the architectures are listed in the Appendix. For all adversarial training, we set batch size to 8 and learning rate to $5e-5$. After obtaining these trained generators, we use them to train a ResNet34 classifier with adversarial training.

Examples of simulated adversarial examples In Figure 2, we show examples generated by our simulator and simulation attack. We can see that the our adversarial examples are of the similar quality as the random ones.

²<https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix>