
[Re] Diffusion-Based Adversarial Sample Generation for Improved Stealthiness and Controllability

William Kang,
syu@student.ubc.ca

Christina Yang
chryang@student.ubc.ca

Abstract

We are doing a reproducibility report based on the paper "Diffusion-Based Adversarial Sample Generation for Improved Stealthiness and Controllability" by Haotian Xue, Alexandre Araujo, Bin Hu, Yongxin Chen. Their Github Repo is here: <https://github.com/xavihart/Diff-PGD/tree/main>.

The paper uses a novel framework to generate adversarial samples. They use a gradient based method guided by a pre-trained diffusion model to try to generate images that appear realistic to the human eye, can fool a wide range of models, and is easy to control how certain regions are modified.

In particular, we are evaluating the claims that Diff-PGD outperform baseline methods such as PGD, AdvPatch, and AdvCam in physical-world attacks and style-based attacks specifically. Finally, we are evaluating the claim that Diff-PGD generates adversarial samples with higher stealthiness.

The following section formatting is optional, you can also define sections as you deem fit. Focus on what future researchers or practitioners would find useful for reproducing or building upon the paper you choose.

1 Introduction

A few sentences placing the work in high-level context. Limit it to a few paragraphs at most; your report is on reproducing a piece of work, you don't have to motivate that work.

2 Scope of reproducibility

The main claims from the original paper are as follows:

1. Diff-PGD can be applied to specific tasks such as digital attacks, physical-world attacks, and style-based attacks, outperforming baseline methods such as PGD, AdvPatch, and AdvCam.
2. Diff-PGD is more stable and controllable compared to existing methods for generating natural-style adversarial samples.
3. Diff-PGD surpasses the original PGD in Transferability and Purification power
4. Diff-PGD generates adversarial samples with higher stealthiness

We will be exploring the first claim in regards to physical-world and style-based attacks specifically.

29 **3 Methodology**

30 We used the authors' code, however, since not all parts of the authors' code have been released, we
31 wrote our own code to evaluate a lot of the experiments.

32 **3.1 Model descriptions**

33 We used ResNet-50, ResNet-101, ResNet-18, Wide-ResNet-50, and Wide-ResNet-101 to classify
34 images. The parameter 'weights' was set to the default weight for each model.

35 **3.2 Datasets**

36 The original paper used ImageNet, but due to limitations in compute and memory, we used a smaller
37 subset of ImageNet with 1000 samples: [https://www.kaggle.com/datasets/figotini/](https://www.kaggle.com/datasets/figotini/imagenetmini-1000)
38 [imagenetmini-1000](https://www.kaggle.com/datasets/figotini/imagenetmini-1000)

39 **3.3 Hyperparameters**

40 Describe how the hyperparameter values were set. If there was a hyperparameter search done, be
41 sure to include the range of hyperparameters searched over, the method used to search (e.g. manual
42 search, random search, Bayesian optimization, etc.), and the best hyperparameters found. Include the
43 number of total experiments (e.g. hyperparameter trials). You can also include all results from that
44 search (not just the best-found results).

45 **3.4 Experimental setup and code**

46 Include a description of how the experiments were set up that's clear enough a reader could replicate
47 the setup. Include a description of the specific measure used to evaluate the experiments (e.g. accuracy,
48 precision@K, BLEU score, etc.). Provide a link to your code.

49 We ran the code given by the authors. We copied the hyperparameter setup of the authors.
50

51 **3.4.1 Success Attack Rate**

52 We wrote our own code to test success rate in `attack_global.py`. We used the same hyperparamete-
53 rs for running Attack_Global as the authors, except that we chose to use `skip=20`, to use 78 images
54 per iteration to test success rate. If the adversarial sample caused the model to classify the image as
55 anything other than it's original classification, we considered it a successful attack.
56

57 **3.4.2 Physical-World Attacks**

58 We tried our own physical world attack using an image patch of a digitally drawn panda head and a
59 laptop as our target object.
60

61 We use a Galaxy A53 5G to take images from the real world and used a MP C4505ex Color Laser
62 Multifunction Printer to print the images in color.
63

64 We stuck the original image on ..., and classified it using ResNet-50.

65 We tested for Success Attack Rate of Diff-PGD using 250 uniformly sampled images from our
66 dataset. (See figure ...)

67 The code for the figure in the paper was not provided, so we created our own code to generate the
68 figure.
69

70 We also need to generate anti-purification table from paper, but I'm not sure how to generate this.

71 Transferability: Figure 6b+6c
 72 We also test the success rate attacking adversarially trained ResNet-50

73 3.5 Computational requirements

74 Include a description of the hardware used, such as the GPU or CPU the experiments were run on.
 75 For each model, include a measure of the average runtime (e.g. average time to predict labels for a
 76 given validation set with a particular batch size). For each experiment, include the total computational
 77 requirements (e.g. the total GPU hours spent). (Note: you'll likely have to record this as you run
 78 your experiments, so it's better to think about it ahead of time). Generally, consider the perspective of
 79 a reader who wants to use the approach described in the paper — list what they would find useful.

80 4 Results

81 Start with a high-level overview of your results. Do your results support the main claims of the
 82 original paper? Keep this section as factual and precise as possible, reserve your judgement and
 83 discussion points for the next "Discussion" section.

Sample	Original paper results				
	(+P)ResNet50	(+P)ResNet101	(+P)ResNet18	(+P)WRN50	(+P)WRN101
x_{PGD}	0.35	0.18	0.26	0.20	0.17
x_n (Ours)	0.35	0.18	0.26	0.20	0.17
x_n^0 (Ours)	0.35	0.18	0.26	0.20	0.17

86 4.1 Results reproducing original paper

87 For each experiment, say 1) which claim in Section 2 it supports, and 2) if it successfully reproduced
 88 the associated experiment in the original paper. For example, an experiment training and evaluating a
 89 model on a dataset may support a claim that that model outperforms some baseline. Logically group
 90 related results into sections.

91 skip=20

92 4.1.1 Attack success rate

93 We got a very different result for success rate from that of the original paper (fig. 1). In the original
 94 paper, the success rate of all of the attacks reached 100% after 5 iterations. However, in our case the
 95 success rate reached only 60%-80% after 15 iterations.

96 Since the authors have yet to release their code for success rate, we are unsure what caused this
 97 difference. We did use a different dataset than the authors, this could be a possible reason.

99 4.1.2 Physical-World Attacks

100 To verify that Diff-PGD can be applied to physical-world attacks, we tried an attack a laptop as our
 101 target object.

102 We used a laptop as our target object and printed out an image of a laptop. We then stuck the image
 103 on a laptop and classified it using all 5 classifiers. We found that the attack was successful for all 5
 104 classifiers. This supports the claim that Diff-PGD can be applied to physical-world attacks.

105 We present the result of physical world attacks using adversarial patches generated using Diff-PGD
 106 in the main paper. In order to show that the adversarial patches are robust to camera views, we
 107 show more images taken from different camera views in Figure 16. For the two cases: computer
 108 mouse (untargeted) and back bag (targeted to Yorkshire Terrier), we randomly sample ten other
 109 camera views. The results show that the adversarial patches are robust both for targeted settings and
 110 untargeted settings. For targetted settings, our adversarial patch can fool the network to predict back
 111 bag as terriers, and for untargeted settings, the adversarial patch misleads the network to predict
 112 computer mouse as artichoke

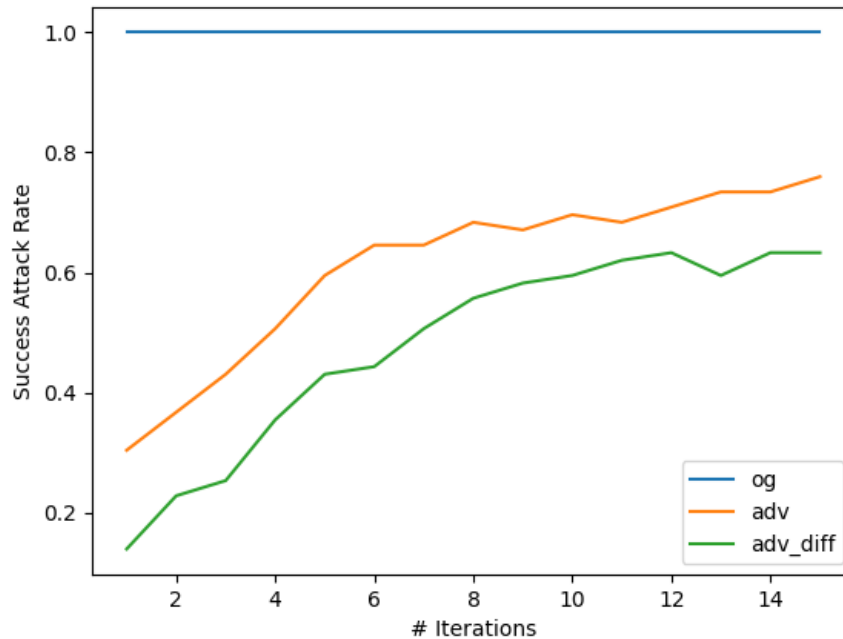


Figure 1: Blue line is success rate of classifying an image not attacked. Orange line is success rate of causing a different classification with Diff-PGD. Green line is success rate of causing a different classification with Diff-PGD after applying sedit at the end (purification??).

4.1.3 Result 2

4.2 Results beyond original paper

Often papers don't include enough information to fully specify their experiments, so some additional experimentation may be necessary. For example, it might be the case that batch size was not specified, and so different batch sizes need to be evaluated to reproduce the original results. Include the results of any additional experiments here. Note: this won't be necessary for all reproductions.

4.2.1 Additional Result 1

4.2.2 Additional Result 2

5 Discussion

Give your judgement on if your experimental results support the claims of the paper. Discuss the strengths and weaknesses of your approach - perhaps you didn't have time to run all the experiments, or perhaps you did additional experiments that further strengthened the claims in the paper.

5.1 What was easy

Give your judgement of what was easy to reproduce. Perhaps the author's code is clearly written and easy to run, so it was easy to verify the majority of original claims. Or, the explanation in the paper was really easy to follow and put into code.

Be careful not to give sweeping generalizations. Something that is easy for you might be difficult to others. Put what was easy in context and explain why it was easy (e.g. code had extensive API documentation and a lot of examples that matched experiments in papers).

132 **5.2 What was difficult**

133 List part of the reproduction study that took more time than you anticipated or you felt were difficult.

134 Be careful to put your discussion in context. For example, don't say "the maths was difficult to
135 follow", say "the math requires advanced knowledge of calculus to follow".

136 The physical-world attack was a lot more difficult to reproduce than we anticipated. AdvCam and
137 AdvPatch both ran fairly quickly, averaging around 2-3s per iteration. However, Diff-PGD averaged
138 around 60s per iteration, which took multiple days to run. Because of this, instead of running 4000
139 iterations as was done in the original paper, we only ran 1500 iterations for each image.

140 None of the models were able to perturb the classification of the image, which was likely due to the
141 fact that we did not adjust the scale properly. This likely means that hyperparameter tuning is very
142 important for physical-world attacks. However, tuning the model is more challenging for Diff-PGD
143 for physical-world attacks, due to how slowly it runs.

144 **5.3 Communication with original authors**

145 Document the extent of (or lack of) communication with the original authors. To make sure the
146 reproducibility report is a fair assessment of the original research we recommend getting in touch
147 with the original authors. You can ask authors specific questions, or if you don't have any questions
148 you can send them the full report to get their feedback before it gets published.

149 **References**

150 **A Supplementary material**

151 **A.1 Physical attacks**