
[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 We ran the code given by the authors. We copied the hyperparameter setup of the authors.
47

48 **3.4.1 Success Attack Rate**

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

54 **3.4.2 Physical-World Attacks**

55 We tried our own physical world attack using an image patch of a digitally drawn panda head and a
56 laptop as our target object.
57

58 We use a Galaxy A53 5G to take images from the real world and used a MP C4505ex Color Laser
59 Multifunction Printer to print the images in color.
60

61 Due to time and compute constraints, we chose to run 1500 iterations for each method instead of 4000
62 like the original model. We chose 1500 because after running one of the original paper's physical
63 world attacks, the loss converged after around 1000 iterations.

64 We took the adversarial samples generated using each method (AdvCam, AdvPatch, and Diff-PDG)
65 and stuck them on a laptop. We then took multiple pictures of the laptop with the sample on top, with
66 the sample in various different locations and rotations, and tested for success rate on these photos
67 by classifying them with ResNet-50.
68

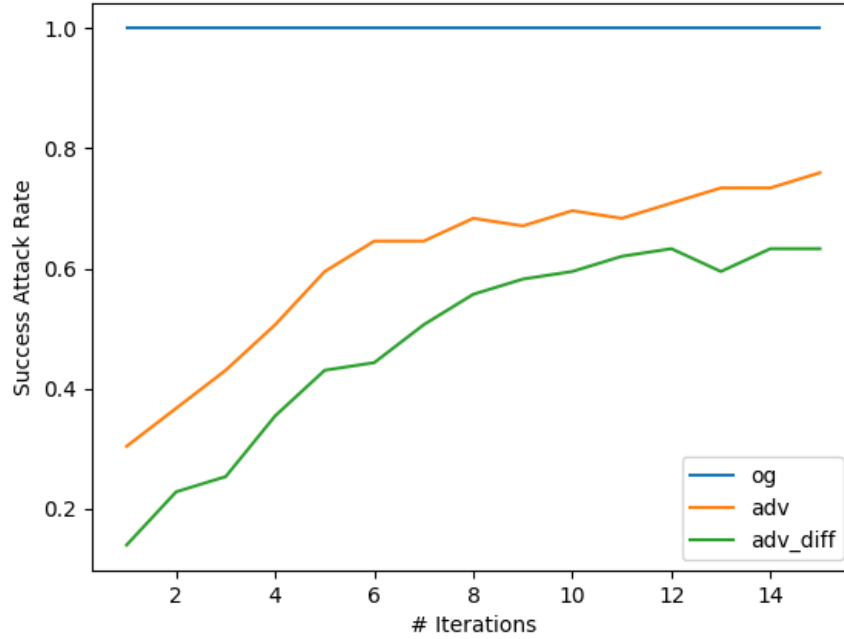


Figure 1: Blue line is success rate of classifying an image not attacked. Orange line is success rate of causing a different classification with Diff-PGD. Green line is success rate of causing a different classification with Diff-PGD after applying sedit at the end (purification??).

69 3.4.3 Style-Based Attacks

70 3.5 Computational requirements

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

77 4 Results

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

81 4.1 Results reproducing original paper

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

86 skip=20

87 **4.1.1 Attack success rate**

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

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

94 **4.1.2 Physical-World Attacks**

95 None of the adversarial samples generated by AdvCam, AdvPatch, or Diff-PGD were able to perturb
96 the classification of the image. This was likely due to the fact that we did not adjust the scale properly.
97 This likely means that hyperparameter tuning is very important for physical-world attacks.

99 **4.1.3 Style-Based Attacks**

100 **4.2 Results beyond original paper**

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

105 **4.2.1 Additional Result 1**

106 **5 Discussion**

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

110 **5.1 What was easy**

111 Running the code from the original authors to generate adversarial samples was fairly easy, since
112 documentation was included in how to run each attack.

114 **5.2 What was difficult**

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

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

118 There was little to no documentation on the inner workings of the code, making it more difficult to
119 understand a lot of the attack parameters, and what they did, however the naming was quite clear for
120 the most part, so this was not too difficult to follow.

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

125 None of the models were able to perturb the classification of the image, which was likely due to
126 the fact that we did not adjust the scale properly. This likely means that hyperparameter tuning
127 is very important for physical-world attacks. However, tuning the model is more challenging for
128 Diff-PGD for physical-world attacks, due to how slowly it runs. Another mistake we made was
129 assuming the the classifier would classify a closed laptop as a laptop, when in reality it classified
130 it as a notebook. We looked into the ImageNet database, and found that most of the images in
131 both laptop and notebook classes were of open laptops, while we used a closed laptop. Since the

132 photo's original label was fed into the model, this could have possibly caused problems with the model.
133

134 **5.3 Communication with original authors**

135 We did not communicate with the original authors.

136 **References**

137 **A Supplementary material**

138 **A.1 Physical attacks**