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

William Kang,
syu@student.ubc.ca

Christina Yang
chryang@student.ubc.ca

Abstract

1 We are doing a reproducibility report based on the paper "Diffusion-Based Ad-
2 versarial Sample Generation for Improved Stealthiness and Controllability" by
3 Haotian Xue, Alexandre Araujo, Bin Hu, Yongxin Chen. Their Github Repo is
4 here: <https://github.com/xavihart/Diff-PGD/tree/main>.
5 The paper uses a novel framework to generate adversarial samples. They use a
6 gradient based method guided by a pre-trained diffusion model to try to generate
7 images that appear realistic to the human eye, can fool a wide range of models, and
8 is easy to control how certain regions are modified.
9 In particular, we are evaluating the claims that Diff-PGD outperform baseline
10 methods such as PGD, AdvPatch, and AdvCam in physical-world attacks and
11 style-based attacks specifically. Finally, we are evaluating the claim that Diff-PGD
12 generates adversarial samples with higher stealthiness.

Reproducibility Summary

13
14 *Template and style guide to ML Reproducibility Challenge 2020. The following section of Repro-*
15 *ducibility Summary is **mandatory**. This summary **must fit** in the first page, no exception will be*
16 *allowed. When submitting your report in OpenReview, copy the entire summary and paste it in the*
17 *abstract input field, where the sections must be separated with a blank line.*

18 Scope of Reproducibility

19 State the main claim(s) of the original paper you are trying to reproduce (typically the main claim(s)
20 of the paper). This is meant to place the work in context, and to tell a reader the objective of the
21 reproduction.

22 Methodology

23 Briefly describe what you did and which resources you used. For example, did you use author's code?
24 Did you re-implement parts of the pipeline? You can also use this space to list the hardware used,
25 and the total budget (e.g. GPU hours) for the experiments.

26 Results

27 Start with your overall conclusion — where did your results reproduce the original paper, and where
28 did your results differ? Be specific and use precise language, e.g. "we reproduced the accuracy to
29 within 1% of reported value, which supports the paper's conclusion that it outperforms the baselines".
30 Getting exactly the same number is in most cases infeasible, so you'll need to use your judgement to
31 decide if your results support the original claim of the paper.

32 **What was easy**

33 Describe which parts of your reproduction study were easy. For example, was it easy to run the
34 author's code, or easy to re-implement their method based on the description in the paper? The goal
35 of this section is to summarize to a reader which parts of the original paper they could easily apply to
36 their problem.

37 **What was difficult**

38 Describe which parts of your reproduction study were difficult or took much more time than you
39 expected. Perhaps the data was not available and you couldn't verify some experiments, or the
40 author's code was broken and had to be debugged first. Or, perhaps some experiments just take too
41 much time/resources to run and you couldn't verify them. The purpose of this section is to indicate
42 to the reader which parts of the original paper are either difficult to re-use, or require a significant
43 amount of work and resources to verify.

44 **Communication with original authors**

45 Briefly describe how much contact you had with the original authors (if any).

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

49 **1 Introduction**

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

52 **2 Scope of reproducibility**

53 Introduce the specific setting or problem addressed in this work, and list the main claims from the
54 original paper. Think of this as writing out the main contributions of the original paper. Each claim
55 should be relatively concise; some papers may not clearly list their claims, and one must formulate
56 them in terms of the presented experiments. (For those familiar, these claims are roughly the scientific
57 hypotheses evaluated in the original work.)

58 A claim should be something that can be supported or rejected by your data. An example is,
59 "Finetuning pretrained BERT on dataset X will have higher accuracy than an LSTM trained with
60 GloVe embeddings." This is concise, and is something that can be supported by experiments. An
61 example of a claim that is too vague, which can't be supported by experiments, is "Contextual
62 embedding models have shown strong performance on a number of tasks. We will run experiments
63 evaluating two types of contextual embedding models on datasets X, Y, and Z."

64 We investigate the main claims from the original paper, which are:

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

71 Each experiment in Section 4 will support (at least) one of these claims, so a reader of your report
72 should be able to separately understand the *claims* and the *evidence* that supports them.

73 **3 Methodology**

74 Explain your approach - did you use the author's code, or did you aim to re-implement the approach
75 from the description in the paper? Summarize the resources (code, documentation, GPUs) that you
76 used.

77 **3.1 Model descriptions**

78 Include a description of each model or algorithm used. Be sure to list the type of model, the number
79 of parameters, and other relevant info (e.g. if it's pretrained).

80 **3.2 Datasets**

81 For each dataset include 1) relevant statistics such as the number of examples and label distributions,
82 2) details of train / dev / test splits, 3) an explanation of any preprocessing done, and 4) a link to
83 download the data (if available).

84 The original paper used ImageNet, but due to limitations in compute and memory, we used a smaller
85 subset of ImageNet with 1000 samples: <https://www.kaggle.com/datasets/figotini/imagenetmini-1000>
86

3.3 Hyperparameters

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

3.4 Experimental setup and code

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

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

Physical-World Attacks:

We first tried the one of the attacks created by the authors, which was a computer-mouse.

We then tried our own physical world attack using an image patch of a ... and a ... as our target object.

We use an (type of phone here, ex. iPhone 8-Plus) to take images from the real world and use an (type of printer here, ex. HP DeskJet-2752) to print the image in color.

We stuck the original image on ..., and classified it using all 5 classifiers (R50, R101, R18, WR50, WR101)

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

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

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

Transferability: Figure 6b+6c

We also test the success rate attacking adversarially trained ResNet-50

3.5 Computational requirements

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

4 Results

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

Original paper results					
Sample	(+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

4.1 Results reproducing original paper

For each experiment, say 1) which claim in Section 2 it supports, and 2) if it successfully reproduced the associated experiment in the original paper. For example, an experiment training and evaluating a

130 model on a dataset may support a claim that that model outperforms some baseline. Logically group
131 related results into sections.

132 **4.1.1 Result 1**

133 **4.1.2 Result 2**

134 **4.2 Results beyond original paper**

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

139 **4.2.1 Additional Result 1**

140 **4.2.2 Additional Result 2**

141 **5 Discussion**

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

145 **5.1 What was easy**

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

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

152 **5.2 What was difficult**

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

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

156 **5.3 Communication with original authors**

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

161 **References**