

---

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

---

William Kang,  
syu@student.ubc.ca

Christina Yang  
chryang@student.ubc.ca

## Reproducibility Summary

1  
2 *Template and style guide to ML Reproducibility Challenge 2020. The following section of Repro-*  
3 *ducibility Summary is **mandatory**. This summary **must fit** in the first page, no exception will be*  
4 *allowed. When submitting your report in OpenReview, copy the entire summary and paste it in the*  
5 *abstract input field, where the sections must be separated with a blank line.*

### 6 Scope of Reproducibility

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

### 10 Methodology

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

### 14 Results

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

### 20 What was easy

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

### 25 What was difficult

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

32 **Communication with original authors**

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

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

## 37 **1 Introduction**

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

## 40 **2 Scope of reproducibility**

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

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

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

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

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

## 61 **3 Methodology**

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

### 65 **3.1 Model descriptions**

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

### 68 **3.2 Datasets**

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

72 The original paper used ImageNet, but due to limitations in compute and memory, we used a smaller  
73 subset of ImageNet with 1000 samples: <https://www.kaggle.com/datasets/figotini/imagenetmini-1000>  
74

### 75 3.3 Hyperparameters

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

### 81 3.4 Experimental setup and code

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

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

86  
87 Physical-World Attacks:

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

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

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

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

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

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

98

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

100 Transferability: Figure 6b+6c

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

### 102 3.5 Computational requirements

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

## 109 4 Results

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

113

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

114

### 115 4.1 Results reproducing original paper

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

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

#### 120 **4.1.1 Result 1**

#### 121 **4.1.2 Result 2**

### 122 **4.2 Results beyond original paper**

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

#### 127 **4.2.1 Additional Result 1**

#### 128 **4.2.2 Additional Result 2**

## 129 **5 Discussion**

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

### 133 **5.1 What was easy**

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

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

### 140 **5.2 What was difficult**

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

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

### 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**