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

William Kang

wkang01@student.ubc.ca

Christina Yang

chryang@student.ubc.ca

Abstract

1 We are doing a reproducibility report based on the paper "Diffusion-Based Adver-
2 sarial Sample Generation for Improved Stealthiness and Controllability" by Haotian
3 Xue, Alexandre Araujo, Bin Hu, Yongxin Chen. The paper uses a novel framework
4 to generate adversarial samples. They use a gradient based method guided by a
5 pre-trained diffusion model to try to generate images that appear realistic to the
6 human eye, can fool a wide range of models, and is easy to control how certain
7 regions are modified. We find that attacks have to be carefully created, with the
8 right hyper-parameters, and target images, in order to get results similar to the
9 original paper.

10 1 Introduction

11 Neural networks are remarkably good at tasks such as image classification. However, they are
12 susceptible to attacks that perturb images in subtle ways, causing the model to classify one image
13 as something else entirely. These perturbed samples are called adversarial samples, and are a major
14 security concern for systems that use neural networks. There are various types of adversarial attacks,
15 including global attacks, which perturb any part of the image, regional attacks, which change a
16 specified region of an image, style-based attacks, which change an image region based on a reference
17 image, and physical world attacks, which generates images you can print out and place on a physical
18 object which cause the physical object to be misclassified by an image classifier.

19 The authors of the original paper propose a new method for generating adversarial models, that uses
20 an off-the-shelf diffusion model to help guide the generation of the perturbations in the image, called
21 Diff-PGD (Xue et al. 2023). This method aims to create adversarial samples that are more realistic,
22 and harder to detect by the human eye, compared to baseline methods, and can be applied to various
23 types of attacks, such as digital attacks, physical-world attacks, and style-based attacks.

24 2 Scope of reproducibility

25 The main claims from the original paper are as follows:

- 26 1. Diff-PGD can be applied to specific tasks such as digital attacks, physical-world attacks,
27 and style-based attacks, outperforming baseline methods such as PGD (Madry et al. 2019),
28 AdvPatch (Brown et al. 2018), and AdvCam (Lee, Kim, and Yoon 2021).
- 29 2. Diff-PGD is more stable and controllable compared to existing methods for generating
30 natural-style adversarial samples.
- 31 3. Diff-PGD surpasses the original PGD in Transferability and Purification power
- 32 4. Diff-PGD generates adversarial samples with higher stealthiness

33 We explore the first claim in regards to physical-world and style-based attacks specifically, and test
34 these attacks on a broader range of images.

35 **3 Methodology**

36 **3.1 Model descriptions**

37 We used ResNet-50 (He et al. 2016) as our primary classifier. We used the pre-trained models with
38 their default weights.

39 **3.2 Datasets**

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

43 **3.3 Experimental setup and code**

44 We used the authors' code (<https://github.com/xavihart/Diff-PGD/tree/main>), however,
45 since not all parts of the authors' code have been released, we wrote our own code to evaluate a lot
46 of the experiments. We also expanded on some of the authors' experiments with our own code, and
47 reduced the number of iterations in a number of cases, due to time and compute constraints. We also
48 added our own code for evaluating the success rates of the authors' code along with creating various
49 visualizations.

50 **3.3.1 Attack Success Rate**

51 We wrote our own code to test success rate. Our measure of success rate is the percentage of
52 successful attacks. We consider an attack successful if the adversarial sample caused the model to
53 classify the image as anything other than its original classification.

54 **3.3.2 Physical-World Attacks**

55 We tried two of our own physical world attacks using an image patch of a digitally drawn panda head
56 and a laptop as our target object for the first one, and an image patch of a forest photo and a water
57 bottle in the second.

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

63 We took the adversarial samples generated using each method (AdvCam, AdvPatch, and Diff-PDG)
64 and stuck them on a the target object. For each adversarial sample (one for each method), we took
65 multiple pictures of the target object with the adversarial sample on top, with the sample in various
66 different locations and rotations, and tested for success rate on the these photos by classifying them
67 with ResNet-50. We also did this for the original unperturbed image patch and the photo with no
68 image patch.

69 **3.3.3 Style-Based and Global Attacks**

70 For each iteration, we take a random base image from the test set of the ImageNetMini dataset. Then
71 for each image we compare the success rate for creating adversarial samples under 3 different attacks:

- 72 1. Style-based attack with a different target image from ImageNetMini. Along with a mask
73 generated based on the base image
- 74 2. The same as the above, except no mask is applied.
- 75 3. Global attack which has no target image and no mask

76 We repeat this for 100 images.

77 **3.4 Computational requirements**

78 We ran it on either CPU or MPS. For used MPS for style-based attacks and global attacks, and CPU
 79 for physical-world attacks.

80 Global attack success rate took between 4-5 hours to run, while physical-world attacks for one target
 81 object and patch took around 1-2 days to run for 1500 iterations.

82 The style-based attack took around 15s per image.

83 **4 Results**

84 **4.1 Physical-World Attacks**

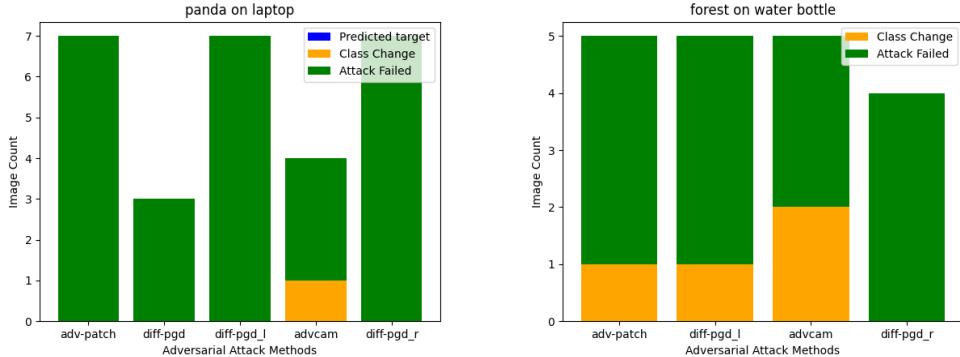


Figure 1: **Evaluating Success Rate of Physical-World Attacks**

Blue line is the number of images that were classified correctly as the target label. Orange line is the number of images that had a different classification with ResNet-50, but that is not the target class. Green line is the number of images whose classifications were one of the following: 'notebook, notebook computer', 'laptop, laptop computer', 'iPod' for panda on laptop, and 'water bottle', 'lighter, light, igniter, ignitor', 'rain barrel', 'pill bottle' and 'hand blower, blow dryer, blow drier, hair dryer, hair drier' for forest on water bottle. These are the classes ResNet-50 gave to images with no adversarial sample (the original image on target object or no image on target object), or classes that we decided had a nearly identical similar meaning as the target object class (ie. laptop vs. notebook computer)

85 For the panda image adversarial sample and the laptop as the target object, at first we found that
 86 none of the adversarial samples generated could perturb the classification of the image. We thought
 87 this might be due to the scale of the image, so we tried to adjust the scale of the image for AdvCam.
 88 From the different camera rotations, zooms and angles we took the AdvCam image from, only one
 89 classified it as something different. However, this class was 'hard disk' (fig. 6), which is also not too
 90 far from the original image's classification of 'notebook, notebook computer'.

91 For the forest image adversarial sample and the water bottle as the target object, we found that the
 92 AdvCam, AdvPatch and Diff-PGD without purification samples were able to perturb the classification
 93 of the image on 1 or 2 of the images, while the Diff-PGD sample with purification did not.

94 **4.2 Style-Based Attacks**

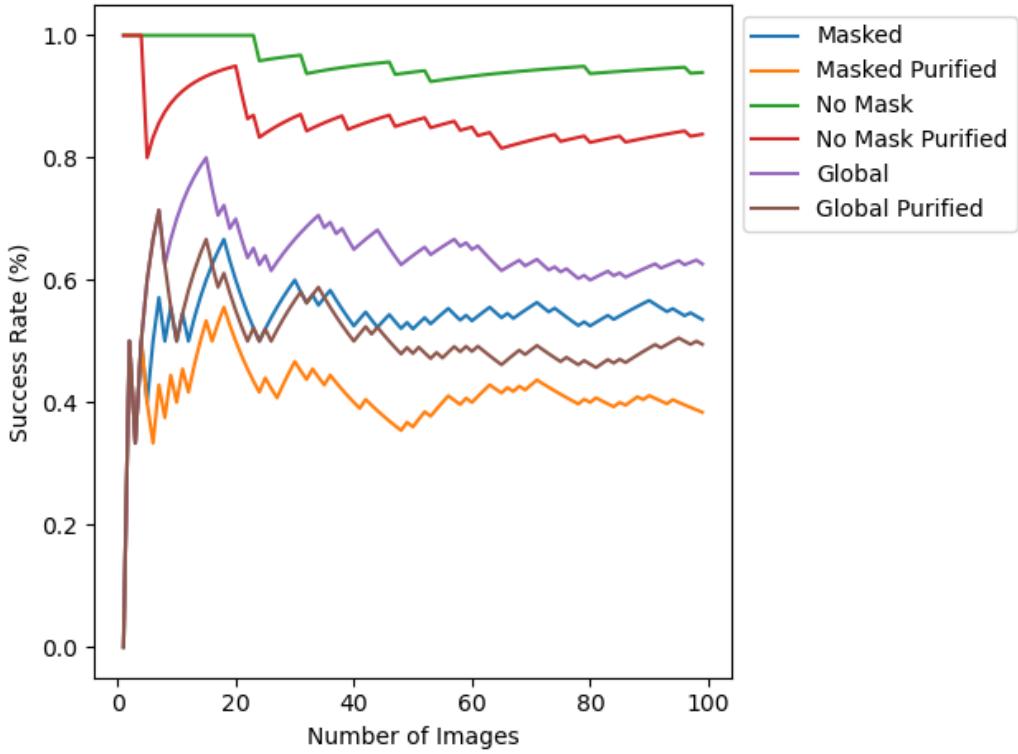


Figure 2: **Plot of the Success Rates:**

We see the best results for applying a style-based attack with no mask. A global attack with no style-image gives the worst results. The success rate for purified images is lower than their non-purified versions.

95 **4.3 Results for Global Attack**

96 From fig. 2 it shows the success rate for global attacks is around 65% before purification, and around
97 50% after purification. This success rate is slightly lower than in the original paper. We used the
98 same hyperparameters as the author, but tested less images.

99 **4.4 Results for Style Attack With Mask**

100 From fig. 2 it shows the success rate for style attacks with a mask applied at around 60% before
101 purification, and around 40% after purification. The original paper does not provide success rates for
102 this type of attack in particular.

103 **4.5 Results for Style Attack Without Mask**

104 From fig. 2 it shows the success rate for style attacks with no mask applied at around 95% before
105 purification, and around 85% after purification. The original paper does not provide success rates for
106 this type of attack in particular. But the anti-purification results hold up as there is minimal drop in
107 success rate.

108 **5 Discussion**

109 **5.1 Attack Success Rates of Physical Attacks**

110 Results were very similar across all three methods (Diff-PGD, AdvCam, AdvPatch), so we failed
111 to support or reject the claim that Diff-PGD can be applied to physical-world attacks. However,
112 from the samples generated, the Diff-PGD samples looked less stealthy and more perturbed than
113 those of AdvCam and AdvPatch, which does not support the paper's claim that Diff-PGD generates
114 samples with higher stealthiness. Fortunately, the stealthiness of physical-world samples is likely
115 less important, since the samples generated through every method were quite obviously perturbed.
116 It also appears that the rotation, zoom, and angle the photo was taken at could affect whether the
117 attack was successful or not, since in the water bottle case, some successful misclassifications were
118 caused, but we did not notice any obvious pattern on what camera angles were better than others.
119 Even without the adversarial samples, the classifications for the water bottle case were quite varied,
120 and often did not make much sense (i.e. 'rain barrel', 'hand blower, blow dryer, blow drier, hair dryer,
121 hair drier'), therefore a likely problem is that we used a dataset that was too small. It may be that
122 since physical attacks have more requirements on robustness, having enough training data is very
123 important on getting the attack to work.

124 **5.2 Attack Success Rates of Global Attack**

125 We did not get great success rates with the global attack. The global attack just adds some noise to
126 the original to change the classification. However, one caveat with this method is that some labels
127 in ImageNet are very similar. For example, there is a difference between "great white shark" and
128 "tiger shark". We observed that even when a global attack was considered "successful", the label
129 would have simply changed to a different label that is still very similar to the original. So while this
130 technically changed the classification to a wrong one, this is an honest mistake that humans would
131 also make.

132 Due to this, a better measure of success might be to compare the distance between the classification
133 labels before and after the attack using a word-embedding model. This way, mistakes such as getting
134 "great white shark" versus "tiger shark" wrong is punished less. Doing this may lead to a lower
135 success score of global attacks since the changes it makes are not very pronounced.

136 **5.3 Attack Success Rates of Style-Based Attacks**

137 Style-Based attacks when successful would generally drastically change the classification label from
138 the original. Attacks with no mask that allowed changes to any part of the image were generally more
139 successful, but greatly loses the "stealthiness" property claimed in the original paper where changes
140 are hard to detect with the human eye. Attacks with a mask applied, are harder to detect but also have
141 a lower success rate.

142 This could be due to the fact that we are using random images as our original and our target. In the
143 original paper, they have specific examples where the target image is of a similar shape to the original.
144 Which leads to the attacks being more effective see fig. 4 for an example.

145 **5.4 What was easy**

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

148 The authors did not include much documentation for how the code worked within the code itself,
149 however they did include some explanations in the paper. This separation made it a bit difficult to
150 understand what the attack parameters did, however the naming was quite clear for the most part, so
151 this was not too difficult to follow.

152 The original code had print statements at each iteration, providing the speed each iteration was taking.
153 This was helpful in timing how long the model would take to run.

154 **5.5 What was difficult**

155 The physical-world attack was a lot more difficult to reproduce than we anticipated. AdvCam and
156 AdvPatch both ran fairly quickly, averaging around 2-3s per iteration. However, Diff-PGD averaged
157 around 60s per iteration, which took multiple days to run. We got very low success rates for the
158 physical-world attacks, which may be because physical-world attacks have more requirements on
159 robustness. This likely means that hyperparameter tuning is very important for physical-world attacks.
160 However, tuning the model is more challenging for Diff-PGD for physical-world attacks, due to how
161 slowly it runs. Another mistake we made was assuming the classifier would classify a closed
162 laptop as a laptop, when in reality it classified it as a notebook. We looked into the ImageNet database,
163 and found that most of the images in both laptop and notebook classes were of open laptops, while
164 we used a closed laptop. Since the photo's original label was fed into the model, this could have
165 possibly caused problems with the model.

166 **6 Conclusions**

167 We failed to reproduce results for physical-world attacks. Our success-rate was very low, and the
168 adversarial samples generated were less stealthy than baseline models. We conclude our hyper-
169 parameter tuning was insufficient to deal with the robustness required for physical-world attacks.

170 For style-based attacks, the Diff-PGD adversarial samples tend to have higher success rates, and are
171 harder to detect with the human eye when the target image is created more-carefully to look like
172 the original. Using random images as the target still lead to success in confusing the model, but
173 oftentimes, the adversarial sample would not look like the original. Also, the measure of success rate
174 could be improved. Some images are hard to classify, such as different breeds of dogs. Perhaps the
175 model should not be punished as heavily for miss-classifying similar labels. This could be done by
176 measuring vector-embedding distance between labels as a measure of success.

177 **7 Acknowledgement**

178 We would like to thank the CPSC 440 teaching team for their guidance this term. We would also like
179 to thank the original authors for their work leading to this paper.

180 **References**

- 181 Brown, Tom B., Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer (2018). *Adversarial
182 Patch*. arXiv: 1712.09665 [cs.CV].
- 183 He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun (June 2016). “Deep Residual Learning
184 for Image Recognition”, pp. 770–778. DOI: 10.1109/CVPR.2016.90.
- 185 Lee, Jungbeom, Eunji Kim, and Sungroh Yoon (2021). *Anti-Adversarially Manipulated Attributions
186 for Weakly and Semi-Supervised Semantic Segmentation*. arXiv: 2103.08896 [cs.CV].
- 187 Madry, Aleksander, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu
188 (2019). *Towards Deep Learning Models Resistant to Adversarial Attacks*. arXiv: 1706.06083
189 [stat.ML].
- 190 Xue, Haotian, Alexandre Araujo, Bin Hu, and Yongxin Chen (2023). “Diffusion-Based Ad-
191 versarial Sample Generation for Improved Stealthiness and Controllability”. *arXiv preprint
192 arXiv:2305.16494*.

193 **A Supplementary material**

194 **B Link to Our Github Repo**

195 <https://github.com/Sokole1/CPSC440Proj>

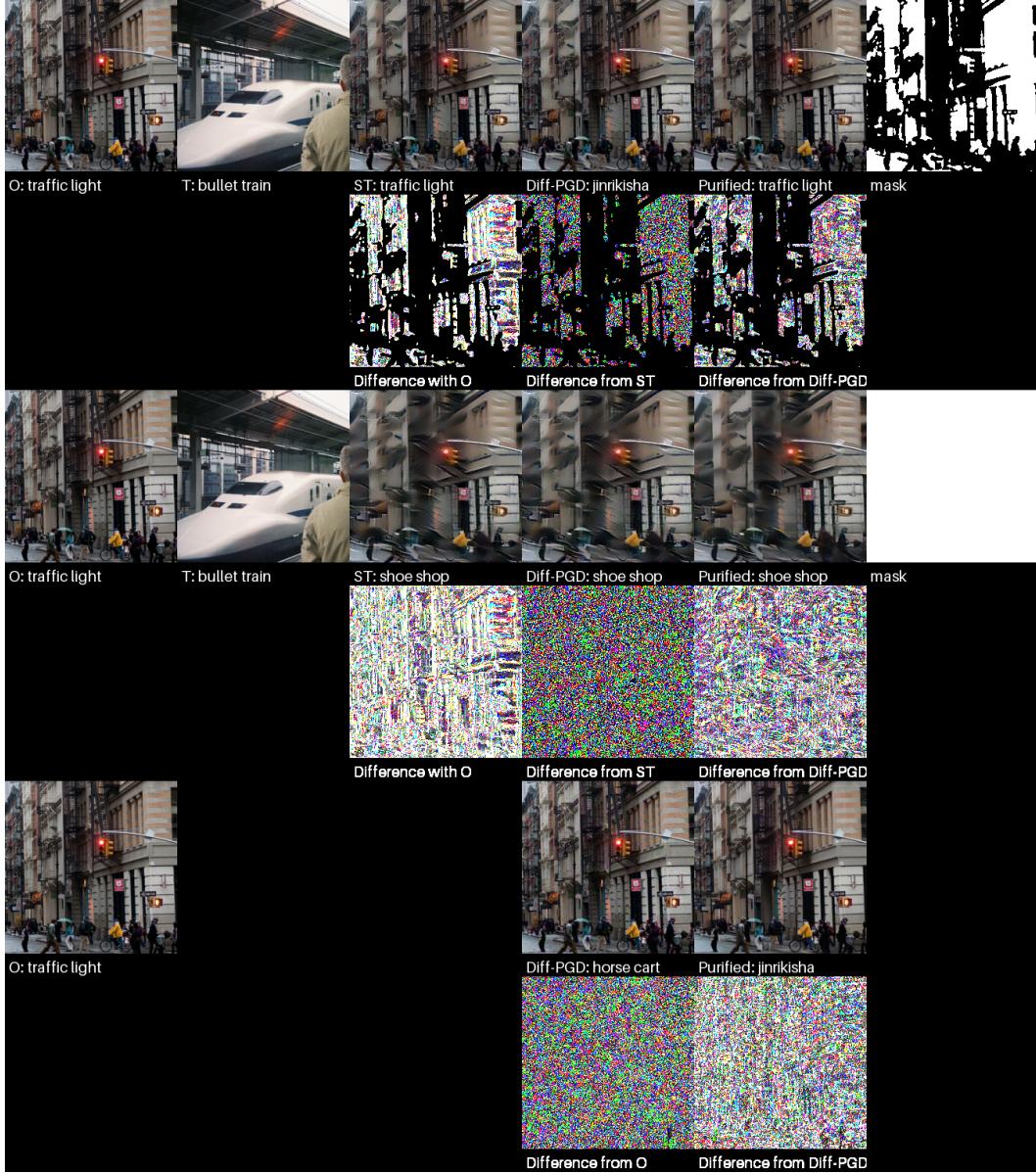


Figure 3: **Visualization of the different attacks:**

There are 6 rows of images. The first two correspond to a style-based attack with a mask applied. The next 2 rows are the same, except there is no mask. The last 2 rows are a global attack where there is no style image. The 6 columns of images are the following: 1: O stands for the original image. 2: T stands for the style image which we try to make the original look more like. 3: ST is the original image with the style image applied, but no adversarial properties. 4: Has the Diff-PGD algorithm applied. 5: Is the image after going through purification. 6: Is the mask applied.

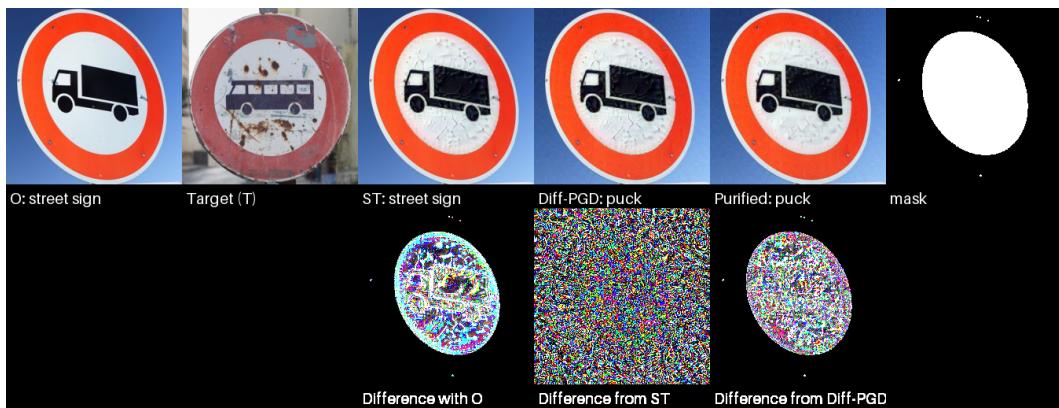


Figure 4: A Hard to Notice Style-Attack

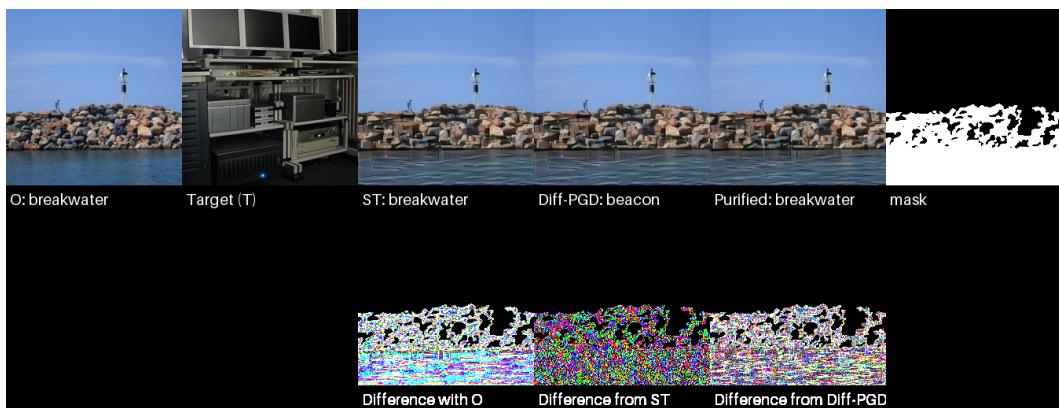


Figure 5: An Easier to Notice Style-Attack

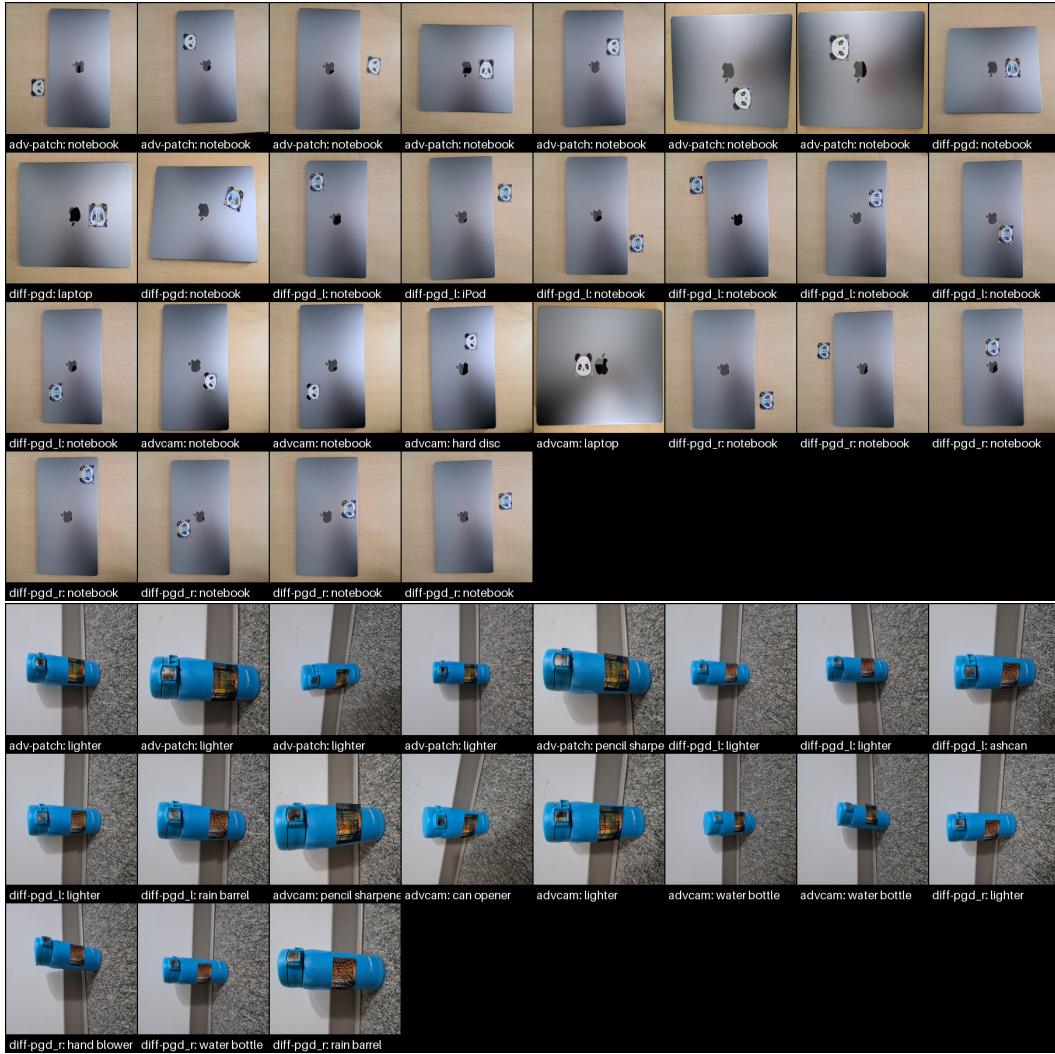


Figure 6: Visualization of the different attacks:
Predictions made for adversarial samples with photos taken at different angles, zooms and rotations.