**1  Rebuttal for 3893: Black-Box Ripper: Copying black-box models using generative evolutionary algorithms**

2  We thank reviewers for their useful comments and insights. The reviewers appreciated our idea as intuitive and original
3  (R1, R3), our method as easy to grasp (R1), significant and highly relevant to NeurIPS (R3), our results as showing the
4  superiority of our method (R1, R3) and our paper as well-written and clear (R1, R2, R3). Furthermore, R3 appreciated
5  our nice review of existing approaches. We next address concerns raised by the reviewers.

6  **R1:** The evolutionary algorithm might collapse to certain regions of the output space. **Answer:** We have analyzed
7  the pixel-level variability of 1000 random images generated by GAN and 1000 random images optimized by our
8  evolutionary strategy, the resulting means of the normalized pixel-level standard deviations being 0.26 and 0.27,
9  respectively. This demonstrates that the proposed evolutionary strategy does not collapse to certain regions.

10  **R1:** Did the authors look at mode collapse of GAN or blurry effect of VAE? **Answer:** The GANs used in our
11  experiments did not suffer from mode collapse. However, for the second set of experiments, the considered VAE
12  does indeed output blurry images. Nonetheless, aspects such as the quality of the images generated by GAN or VAE
13  are not relevant to our approach, as long as we are able to achieve state-of-the-art performance levels in stealing the
14  functionality of black-box models, as shown in Tables 1, 2, 3 from our paper.

15  **R2:** The experimental parts are weak, only small data sets being employed. **Answer:** First, we considered the same
16  data sets as Addepalli et al. [1], our main competitor, their paper presenting results on such data sets being published at
17  AAAI 2020. Aside from the experimental setup used by Addepalli et al. [1], in which the images have low resolution,
18  we have performed additional experiments using larger images from CelebA-HQ, ImageNet Cats and Dogs and 10
19  Monkey Species. In the experiments presented in Table 3 (see paper), the teacher and student models take images of
20  $224 \times 224$ pixels as input. We thus believe that our experiments cover a wider range than those of Addepalli et al. [1].

21  **R2:** Concerns about the proposed method in practical applications with many classes. **Answer:** The experiments
22  performed on higher-resolution image shows that our method works well when confronted with a large multi-class
23  latent space, considering that ImageNet Cats and Dogs contains 143 classes. Regarding our results on 90 and 40 classes
24  from CIFAR-100, we note that our accuracy rates are comparable to those reported by Addepalli et al. [1]. However,
25  our method is significantly more general, as their method only works in white-box scenarios, requiring complete access
26  to the network to perform back-prop. Our method works in a black-box setup, requiring no knowledge of the internal
27  structure or parameters of the model.

28  **R2:** Is the GAN used in the paper same as in the baseline [1]? **Answer:** The GANs have equivalent architectures, but
29  the models are not identical, mainly because Addepalli et al. [1] train the GAN along with the student, whereas we only
30  use a pre-trained GAN.

31  **R3:** The main weakness is the lack statistics over several independent runs. **Answer:** We note that Addepalli et al. [1]
32  typically reported results of independent runs, with one exception. They performed 5 runs for a single scenario, in
33  which 10 classes from CIFAR-100 are used as proxy. We hereby present results for 5 runs of our approach for the same
34  case in Table 1, demonstrating the stability and superiority of our approach over that of Addepalli et al. [1]. Upon
35  acceptance, we will include mean and standard deviations over 5 runs for all our experiments in the final paper.

Table 1: Accuracy rates (in $\%$) for 5 runs on CIFAR-10 as true data set and 10 classes of CIFAR-100 as proxy.

| Method | Run 1 | Run 2 | Run 3 | Run 4 | Run 5 | Mean $\pm$ Std. Dev. |
|---|---|---|---|---|---|---|
| DeGAN [1] | 66.9 | 74.6 | 72.6 | 76.6 | 71.6 | $72.6 \pm 3.26$ |
| Black-Box Ripper (Ours) | 78.2 | 78.0 | 77.4 | 77.9 | 78.2 | $77.9 \pm 0.29$ |

36  **R3:** Wouldn't it help to perform adaptive mutations? **Answer:** In a set of preliminary experiments, we tried using
37  CMA-ES, without observing any performance improvements. Hence, we decided to stick to the most straight forward
38  method that already outperforms existing approaches.

39  **R3:** Authors only evolve for 10 generations. **Answer:** In a set of preliminary experiments, we tested with up to 50
40  generations, but we observed that the evolved exemplars typically converge in 10 generations or even less. To reduce
41  the computational time, we decided to keep only 10 generations, which seems sufficient (see Tables 1, 2, 3 in paper).

42  **R5:** The detailed description of the evolutionary algorithm used in this paper is not clear, such as what kind of coding
43  method is used. **Answer:** The coding method is straight forward: our data samples are input noise vectors from the
44  latent space of the generator. Upon acceptance, we will strive to further improve the clarity of our presentation.

45  **R5:** Why do you choose this kind of evolutionary algorithm? **Answer:** We tried the simplest evolutionary algorithm,
46  which already seems capable of surpassing glass-box approaches. Further tweaking the evolutionary algorithm can only
47  bring improvements, putting an even greater gap between our results and those of the state-of-the-art methods.

48  **R5:** Paper is not innovative enough. **Answer:** Model stealing has never been studied or proven possible with
49  evolutionary algorithms. We show that our evolutionary strategy surpasses other, more relaxed, glass-box state-of-the-
50  art methods. Furthermore, the originality of our idea is appreciated by R1, R3.

51  [1] Addepalli et al., DeGAN: Data-Enriching GAN for Retrieving Representative Samples from a Trained Classifier.
52  AAAI 2020.