

1 First, we thank all the reviewers for their invaluable assessment of our paper in this challenging time. As they agree, the  
 2 general idea of AdvFlows is sound and promising, and the paper is well-written and self-contained. In the following,  
 3 we address some of the questions raised by the reviewers as much as time and space allows.

4 **Overview** The final goal of designing adversarial attacks is gaining a better insight into the pitfalls of DNNs, ultimately  
 5 alleviating such threats. In this regard, designing attacks with a statistical flavor is extremely valuable as they: 1) provide  
 6 a unifying framework of modeling DNNs’ adversarial vulnerability, and more importantly, 2) help in establishing  
 7 the required connection with mature fields like high-dimensional statistics to use their results in finding the ultimate  
 8 solution to making DNNs more robust. Having these in mind, we have come up with AdvFlow that can be viewed as an  
 9 important step in this direction.

10 **Why NFs and not GANs?** The ability of *Normalizing flows* (NF) for efficient inference and sampling, as well as their  
 11 straightforward and stable training, made them an ideal candidate for our purpose of designing a black-box attack with  
 12 a statistical perspective. Note that *generative adversarial networks* (GAN) have many disadvantages for use in the  
 13 current framework: 1) It is known that GANs suffer from *mode collapse*, where they fail to represent different modes of  
 14 data equally well. In contrast, flow-based models are trained to maximize the log-likelihood, and as such, they cover  
 15 different modes of data better. 2) Finding the latent space representation of data in GANs requires solving a non-convex  
 16 optimization problem by back-propagating through the model for every new attack. However, the proposed NF models  
 17 are invertible by design, and to find the latent space representation of an image, one only needs to query the model.  
 18 3) More importantly, GANs neither represent an explicit distribution nor enable inference and density computation. The  
 19 current design, however, enables further investigation of the attacker distribution properties in the future.

20 **Attack strength** The primary purpose of the current work is to convey the idea of blending statistical methods like  
 21 normalizing flows and adversarial attacks so that we can better understand such threats. Thus, we aimed to compare with  
 22 recent, but widely recognized black-box attacks for comparison. Nevertheless, by doing more rigorous hyper-parameter  
 23 tuning or adding extra variables (like  $\sigma$  as correctly indicated by R2), the results can be improved further.<sup>1</sup>

24 **Ablation study on adversarial example detectors** To provide more reliable evidence that AdvFlow’s distributional  
 25 properties are fooling the adversarial example detectors, we perform the following ablation study. First, we use an  
 26 untrained (denoted by un.) AdvFlow model that is initialized randomly. Then, we use the trained version (denoted by  
 27 tr.) of the same architecture to perform black-box attacks. Using examples generated by these two models, we then  
 28 train adversarial example detectors to spot the adversaries from clean images. In the paper, we used the Mahalanobis  
 29 detector [26], a well-known SOTA adversarial example detector. For the sake of completeness, we also add LID [31]  
 30 (the previous SOTA) and Res-Flow [58] (the recently introduced SOTA) alongside Mahalanobis detector. We compare  
 31 our results with  $\mathcal{N}$ ATTACK, which also approaches the black-box adversarial attack from a distributional perspective  
 32 for a fair comparison. The results are given in Table 1. As shown, only if we pre-train our method on clean data, we can  
 33 fool the detectors. This is indicating that the attacker’s distributional properties are fooling the detectors.

34 **Performance comparison with SimBA [59]** Note that SimBA [59] was not included in the original manuscript as it  
 35 is designed for efficient  $\ell_2$  attacks. At the time of writing the paper, it was not clear how it can be generalized to  $\ell_\infty$ .  
 36 Not until after the NeurIPS deadline did the authors include a generalized version for  $\ell_\infty$ , alongside the explanations.<sup>2</sup>  
 37 We repeat the CIFAR-10 experiments of the paper using the recent version of SimBA-DCT, and report the results in  
 38 Table 2. For a fair comparison, we compute the average and median of queries on examples where both methods have  
 succeeded. As seen, we get similar results to Table 2 of the paper, outperforming SimBA in defended baselines.<sup>3</sup>

Table 1: Adv. example detection on CIFAR-10.

Detector	AUROC(%) $\uparrow$		
	$\mathcal{N}$ ATTACK	Ours (un.)	Ours (tr.)
LID [31]	78.69	84.39	<b>57.99</b>
Mah. [26]	97.95	99.50	<b>66.85</b>
Res. [58]	97.90	99.40	<b>67.03</b>

Table 2: Performance comparison with SimBA [59] on CIFAR-10.

Defense	Success Rate(%) $\uparrow$		Query Avg. $\downarrow$		Query Med. $\downarrow$	
	SimBA	Ours	SimBA	Ours	SimBA	Ours
Vanilla	<b>99.98</b>	<b>99.42</b>	<b>238.08</b>	949.55	<b>126</b>	400
FreeAdv	35.52	<b>41.21</b>	497.97	<b>458.35</b>	256	<b>200</b>
FastAdv	35.07	<b>40.22</b>	<b>469.15</b>	<b>477.77</b>	245	<b>200</b>
RotNetAdv	35.63	<b>40.67</b>	499.75	<b>453.26</b>	267	<b>200</b>

<sup>1</sup>Note that some of the current SOTA results in black-box adversarial attacks come from the attacker’s knowledge about the gradients of the target classifier using substitute models. However, once the target changes its training procedure (e.g., from vanilla to adversarial training), the performance of such methods drop significantly. In contrast, **our method is trained only on clean data and does not depend on any substitute network**. As such, it has a considerable advantage against these methods that are currently prevalent.

<sup>2</sup>See the official repo. of SimBA, where it clearly is indicated that the  $\ell_\infty$  attack is added on 2020/06/22, after NeurIPS deadline.

<sup>3</sup>The results of Table 1 and 2 (as well as SVHN) will be added to the camera-ready version.

[58] Zisselman and Tamar. “Deep Residual Flow for Out-of-Distribution Detection.” *CVPR*, 2020.

[59] Guo et al. “Simple Black-box Adversarial Attacks.” *ICML*, 2019.