

1 We appreciate all reviewers for their helpful and constructive comments. We’ll further improve the paper in the final  
2 version. Below we address their detailed comments.

3 **R1: RGF outperforms NES:** The major difference between RGF and NES [16] is that NES adopts the antithetic  
4 sampling, while RGF does not. Specifically, the gradient estimator is  $\hat{g} = \frac{1}{q} \sum_{i=1}^q \frac{f(x+\sigma u_i, y) - f(x-\sigma u_i, y)}{2\sigma} u_i$  in NES and  
5  $\hat{g} = \frac{1}{q} \sum_{i=1}^q \frac{f(x+\sigma u_i, y) - f(x, y)}{\sigma} u_i$  in RGF (see Eq.(5)). The NES estimator can eliminate the second-order component  
6 of  $f$  through central differences, but it requires  $2q$  queries while RGF only requires  $q + 1$  queries. When  $\sigma$  is small, the  
7 second-order component is often dominated by the first-order one. So RGF outperforms NES. We’ll make it clearer.

8 **R1:  $\lambda^*$  distribution and cosine similarity across the attack iterations:** Thanks  
9 for the suggestion. As it’s hard to plot  
10 the full distribution of  $\lambda^*$ , which changes  
11 during iteration, we show the average  $\lambda^*$   
12 over all images w.r.t. iterations in Fig. A.  
13 It shows that  $\lambda^*$  decreases along with the  
14 iterations (i.e., the distribution concen-  
15 trates on small  $\lambda^*$ ). Fig. B shows the cosine  
16 similarity between the transfer and the true  
17 gradients, across iterations. The results show  
18 that the transfer gradient is useful at begin-  
19 ning, and becomes less useful along with the  
20 iterations. However, the estimated gradient  
21 can remain higher cosine similarity with  
22 the true gradient, which facilitates the adver-  
23 sarial attacks consequently. We’ll add the  
24 results in the final version.

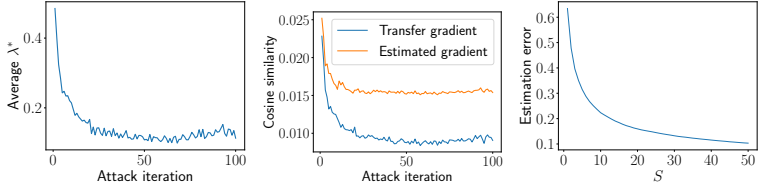


Figure A: The average  $\lambda^*$  across attack iterations. Figure B: The cosine similarity between the transfer and the true gradients, across attack iterations. Figure C: The estimation error with different  $S$ .

16 becomes less useful along with the iterations. However, the estimated gradient can remain higher cosine similarity with the true gradient, which facilitates the adversarial attacks consequently. We’ll add the results in the final version.

20 **R2: Novelty of the idea:** As stated in L108-113, we consider the score-based setting while [4] focuses on the decision-based setting. [4] is built upon the Boundary method [3] and uses a fixed coefficient to incorporate the transfer gradient. Due to the different settings, we introduce a new objective (see Eq.(7)) for gradient estimation, and optimize it inside the proposed family of estimators, resulting in a generic P-RGF algorithm which incorporates the transfer gradient with an optimal coefficient. Technically, it’s non-trivial to derive the optimal solution. Moreover, we found that it’s necessary to use an adaptive coefficient rather than a fixed value since 1) the usefulness of the transfer gradient varies across iterations; 2) experiments show that our algorithm is beneficial from the adaptive coefficient. Overall, we propose a simple, yet novel and effective method, considering a different black-box setting from [4], as agreed by R1 and R3.

28 **R2: More analysis and experiments about the estimation of gradient norm:** Thanks for the comment. The gradient norm (or cosine similarity) is easier to estimate than the true gradient since it’s a scalar value. Fig. C shows

Table A: Additional experimental results.

Methods	Inception-v3		VGG-16		ResNet-50	
	ASR	AVG. Q	ASR	AVG. Q	ASR	AVG. Q
P-RGF ( $\lambda = 0.05$ )	97.8%	1021	99.7%	888	99.6%	790
P-RGF ( $\lambda^*$ , true norm)	98.1%	768	99.8%	501	99.5%	427
P-RGF ( $\lambda^*$ )	98.1%	745	99.8%	521	99.6%	452

32 the estimation error of the gradient norm, defined as the (normalized) RMSE— $\sqrt{\mathbb{E}(\frac{\|\nabla f(x)\|_2 - \|\nabla f(x)\|_2}{\|\nabla f(x)\|_2})^2}$ , w.r.t. the  
33 number of queries  $S$ . We chose  $S = 10$  in all experiments to reduce the number of queries while the estimation error is  
34 acceptable. We also show the overall attack results of using the true gradient norm instead of the estimated norm in  
35 Table A (Row 2). The results are similar to those of using the estimated norm. We’ll add the results in the final version.

36 **R2: Experiments about P-RGF with a fixed  $\lambda = 0.05$ :** Thanks for the suggestion. Table A (Row 1) shows the results  
37 of P-RGF with  $\lambda = 0.05$  (optimal in Fig. 1(b)), which are better than P-RGF with  $\lambda = 0.5$  (in Table 1). However, a  
38 significant performance gap still remains from using the adaptive  $\lambda^*$ . We’ll add the results in the final version.

39 **R3: The improvement over the RGF method is not significant:** In Table A, P-RGF and RGF obtain similar attack  
40 success rates. The reason is that the maximum number of queries (i.e., 10,000) is sufficient for them to find adversarial  
41 perturbations, such that their attack success rates are similarly high. However, P-RGF requires fewer queries than RGF  
42 (20% ~ 45% queries reduction). If the maximum number of queries is set to 1,000, the attack success rate against  
43 Inception-v3 becomes 56.4% using RGF and 78.6% using P-RGF (the average number of queries is 470 and 297  
44 respectively). Moreover, in Table 2, P-RGF obtains much higher success rates than RGF, and also reduces the query  
45 complexity for attacking the defensive models. In summary, the improvement is significant in most of the cases.

46 **R3: Attack results on adversarially trained defensive models:** Thanks for the sug-  
47 gession. We choose [\*1] as our target model, which successfully performs PGD-based  
48 adversarial training on ImageNet. The gradient from ResNet-152 can hardly transfer to  
49 this model, and the results of RGF and P-RGF are similar. So we use another adversarially  
50 trained model (with a different architecture) to provide the transfer gradient. We perform  $\ell_\infty$  attacks with  $\epsilon = 16/255$ ,  
51 which is the same threat model used in adversarial training. Table B presents the results—P-RGF outperforms RGF  
52 significantly with the strong transfer-based prior. We’ll add the results in the final version.

Table B: Attack results on adversarially trained model.

	ASR	AVG. Q
RGF	31.7%	1207
P-RGF ( $\lambda^*$ )	64.7%	378

53 [\*1] C. Xie, Y. Wu, L. van der Maaten, et al. Feature denoising for improving adversarial robustness. CVPR 2019.