

1 We thank the reviewers for their helpful comments. Below, we address some of the points made regarding our work’s  
2 contributions and relation to prior work.

3 **On automated attacks (Reviewer 1 and 3).** Reviewer 1 and Reviewer 3 argue that “AutoAttack” (Croce & Hein,  
4 2020) is a counter-example to our claim that adaptive attacks are inherently necessary and important to study. However,  
5 we believe that the results from that work exactly *support* our claim:

6 1. **AutoAttack fails to attack many defenses we break completely.** AutoAttack does not fully break the challenging  
7 “k-winners take all” defense (19% accuracy), whereas we reduce it to 0% accuracy. Our adaptive attack also  
8 outperforms AutoAttack on “ME-Net” and on at least two other defenses we independently evaluated (“Asymmetrical  
9 Adversarial Training” and “Are Generative Classifiers More Robust”).

10 2. **AutoAttack only applies to standard feed-forward classifiers.** Many defenses depart from the standard defense  
11 template that AutoAttack supports. Of the 13 defenses we study, 5 aim at detecting adversarial examples. For these,  
12 formulating an appropriate loss function to optimize is precisely the challenge in developing a strong attack, and it  
13 is unclear how this process could be automated by AutoAttack.

14 AutoAttack also cannot be directly applied to “Temporal Dependency” (a speech-to-text model) and “Robust Sparse  
15 Fourier Transform” (which is aimed at perturbations of small  $\ell_0$  norm).

16 In fact, a majority of defenses evaluated by AutoAttack use adversarial training, for which there is usually no need  
17 for an “adaptive” attack as the inference phase is unchanged. This addresses a question of Reviewer 1, on why we  
18 refrained from evaluating such defenses in our paper.

19 3. **It is easy to build an ineffective defense that AutoAttack fails to break.** First, pick any of the defenses evaluated  
20 on AutoAttack for which the query-based attacks (FAB and Square) fail. Then, add a component that masks  
21 gradients (e.g., a non-differentiable quantization layer) to defeat PGD. Now AutoAttack fails, even though the  
22 defense is easy to circumvent with an adaptive attack.

23 (This example also addresses another question of Reviewer 1: since query-based attacks routinely fail, e.g., for  
24 randomized defenses, we cannot use them to reliably evaluate all defenses.)

25 We believe AutoAttack is a strong, *non-adaptive* baseline. However, it is orthogonal to our paper, which argues that  
26 static automated attack strategies are not sufficient. The above points illustrate why.

27 **On adversarial models (Reviewer 2).** We apologize for not clarifying this in the paper. All 13 defenses assume a  
28 white-box adversary with full access to the defense parameters. The defenses differ slightly in the perturbation norms  
29 and bounds that they consider, and these are mostly incomparable. As a result, all our attacks assume white-box model  
30 access and operate with the same perturbation bounds as in the original evaluation. While some of our attacks use  
31 black-box optimization techniques, this is solely to side-step gradient-masking. We still view these as white-box attacks.

32 **On related work & technical novelty (Reviewer 3).** We view the fact that “defenses are broken by existing tech-  
33 niques” as a strength rather than a weakness. Introducing new attack techniques would give the incorrect impression  
34 that we currently lack the appropriate technical tools to evaluate defenses properly. But for all of the defenses we  
35 studied, a better attack can be built using only tools that are well-known in the literature.

36 Thus, **the issue with current defense evaluations is methodological rather than technical.** Proposing new attack  
37 techniques is not going to fix this. Instead, our paper clearly documents how to make use of existing techniques to build  
38 strong adaptive attacks.

39 This is what differentiates our work from prior work that proposed and argued for adaptive attacks (e.g., Carlini &  
40 Wagner 2017, Athalye et al. 2018, Carlini et al. 2019). When those papers were written, defense authors indeed lacked  
41 the motivation and tools to conduct strong adaptive evaluations. But of the 13 defenses we study, nearly all claim to  
42 perform an adaptive attack evaluation following current best practices. The main contribution of our paper is thus to  
43 highlight that there still is a systemic lack of strong adaptive evaluations in this field, and to show how to remedy this.