

1 We thank all reviewers for their positive reception of our paper and for their constructive feedback.

2 Reviewer #2

- 3 • **On dual norms and prior work.** Thank you for pointing us to the relevant prior work of Demontis et al. and
4 Xu et al. which we apparently missed.

5 Our results about the tradeoff between ℓ_1 and ℓ_∞ robustness do indeed share similarities with the dual-norm
6 behavior studied by Demontis et al. for linear classifiers. The main difference is that for our Gaussian dataset,
7 we show that the robustness tradeoff is inherent to *any* classifier (i.e., not only linear ones). Another difference
8 is that we extend our results to tradeoffs between ℓ_∞ and spatial perturbations.

9 For arbitrary dual norms ℓ_p and ℓ_q , we can also exhibit a tradeoff as follows (somewhat informally):

- 10 – Flipping the features x_1, \dots, x_n requires a ℓ_p -perturbation of magnitude $O(\mu \cdot n^{1/p}) = O(n^{1/p-1/2})$.
- 11 – Flipping the feature x_0 requires a ℓ_q -perturbation of magnitude $O(1)$.
- 12 – So if a model is robust to perturbations of size $O(n^{1/p-1/2})$ in the ℓ_p -norm, it cannot also be robust to
13 perturbations of size $O(1)$ in the dual norm.

14 We will discuss these connections between our work and the prior work of Demontis et al. and Xu et al. in the
15 final version of our paper. We thank the reviewer for suggesting this dual-norm view which nicely generalizes
16 one of our results.

- 17 • **On structure and readability.** We agree that our paper contains many contributions (the formal analysis, a
18 new ℓ_1 -attack, an experimental evaluation) that are somewhat heterogenous. We will make an effort to clarify
19 our main contributions and to better structure our paper to improve its readability.

20 Reviewer #3

- 21 • **On MNIST artifacts.** The gradient masking effect we discover and explain is indeed specific to MNIST (for
22 multiple ℓ_p norms), and we do not claim otherwise. In fact, this gradient masking effect seems inherently due
23 to the mostly binary nature of the MNIST images, which leads to thresholding being a viable defense against
24 ℓ_∞ -perturbations.

25 Nevertheless, as MNIST is the only vision dataset for which we've been able to train models to high levels of
26 robustness (for individual attack models), we believe it is worthwhile to observe that extending this robustness
27 to multiple ℓ_p attacks may be particularly challenging for this dataset. In this sense, even a dataset as simple as
28 MNIST is clearly not solved from an adversarial robustness perspective.

29 We do believe that it should be possible to train MNIST models to a robustness tradeoff similar to that we
30 found on CIFAR10. But this will require new techniques that somehow circumvent gradient masking as a
31 spurious solution. We think this is an interesting open problem for the community to consider.

32 On CIFAR10, the robustness tradeoff is indeed smaller but still quite noticeable. It is worth noting that our
33 experiments on CIFAR10 only ever consider the combination of two attack types. It is not clear what would
34 happen if we were to try to train models to be robust to 4 or 5 attacks at a time for instance. For these types of
35 experiments we are mainly limited by the poor scalability of adversarial training (e.g., for two attack types,
36 adversarially training a wide ResNet takes about two GPU weeks). There is some promising recent research
37 on speeding up adversarial training, so these types of experiments might become tractable for future work.

- 38 • **On black-box attacks.** The Ensemble Adversarial Training technique of Tramer et al. was proposed to
39 increase a model's robustness to black-box attacks, but it was found to have no noticeable effect on the model's
40 robustness to stronger white-box attacks. As our evaluation focuses on the white-box robustness of the trained
41 models, we have not incorporated black-box attack examples at training time.

42 We also considered using black-box attacks at evaluation time (e.g., as a test against gradient-masking), but
43 found decision-based attacks to be stronger and more reliable for this purpose.