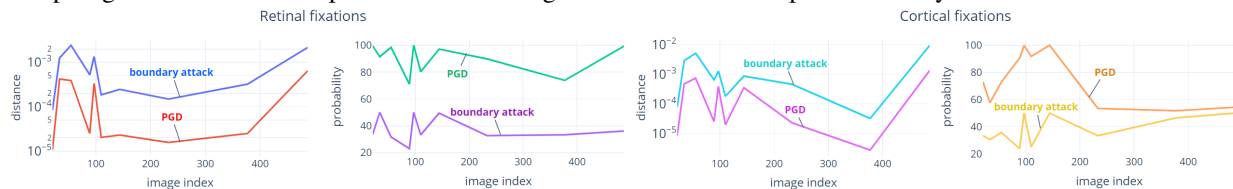1  We thank the reviewers for their insightful and positive feedback. We are very encouraged they enjoyed (R2, R4) the
2  idea to take inspiration from the human visual system to increase robustness of DNNs and found it to be unique (R3) and
3  worthy of investigating more. The reviewers appreciated the substantial amount of control experiments and following
4  best practices better than many other papers (R4), including showing scaling between datasets of varying complexities
5  (R1). We are pleased R2 found the ablation studies informative. We appreciate the constructive comments and questions
6  from the reviewers, some on the organization of the text between the main part and the extensive (R4) Supplementary
7  Material, and some suggesting further experimental tests to strengthen the results. We especially would like to thank
8  R4 for suggesting several very relavant papers to cite. As we were constrained by limited time and computational
9  resources, as well as the space to answer, we focus below on the most crucial additional experiments and questions, but
10 will incorporate all feedback in the final version.

11 **(R2, R4) Run PGD with many more steps and smaller step size to truly evaluate robustness.** The reviewers raised
12 concerns that our main tests of robustness were against the PGD attack with 5 or 20 steps with a large step size of
13 $\epsilon/3$ (we would like to point out that we have verified robustness to smaller step size of $\epsilon/12$ in the Supp. Mat.). Both
14 reviewers suggested to run attacks with smaller step size and follow (Brendel et al, NeurIPS 2019), where attacks were
15 performed up to a 1000 steps. We have followed this suggestion and performed PGD attacks on the ImageNet10 dataset
16 with step size of $\epsilon/20$ (as suggested by R4) and over a range of steps up to 1000, as shown in the Table below. The
17 drop in robustness is minimal with the additional steps. The trend of robustness seen with the small number of steps
18 was already strongly representative of the general trend, and we believe these additional results help drive the message
19 home that both retinal and cortical sampling are biologically-motivated methods that significantly improve robustness at
20 small perturbations without noticeable computational overhead.

Table 1: PGD with more steps on ImageNet10. Standard: standard ResNet; Coarse: coarse fixations on standard
ResNet; Retinal: Retinal fixations model; Cortical: cortical fixations model (Due to time constraints, we were unable
to finish running the 1000 step attack at $\epsilon = 0.001$, but up to 500 steps the accuracy stayed at 82%, with matched
performance for bigger $\epsilon$).

| $\epsilon$ | STEPS = 20, STEP SIZE = $\epsilon/3$ | | | | STEPS = 300, STEP SIZE = $\epsilon/20$ | | | | STEPS = 1000, STEP SIZE = $\epsilon/20$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | STANDARD | COARSE | RETINAL | CORTICAL | STANDARD | COARSE | RETINAL | CORTICAL | STANDARD | COARSE | RETINAL | CORTICAL |
| 0.02 | 0% | 0% | 0.4% | 1% | 0% | 0% | 0.2% | 0.8% | 0% | 0% | 0.2% | 0.8% |
| 0.01 | 1.6% | 1.2% | 7.6% | 11.8% | 1.4% | 0.8% | 7% | 9.4% | 1.4% | 0.8% | 6.8% | 9.4% |
| 0.005 | 20% | 16.4% | 43.4% | 42.8% | 19.4% | 15.8% | 43.4% | 42.2% | 19.2% | 15.4% | 43.4% | 41.8% |
| 0.001 | 78.4% | 80.8% | 86.4% | 82% | 78.4% | 80.8% | 86.4% | 82% | 78.4% | 80.8% | 86.4% | x% |

21 **(R4) Run boundary attack on a larger subset of images and with many more steps.** The reviewer suggested that
22 evaluating the boundary attack on a single image for 1000 iterations is insufficient to claim avoiding gradient obfuscation.
23 While we would claim that it served it's purpose together with many other tests reported in the main text and shown
24 in Supp Mat., we have nonetheless ran the attack on more images for up to 10000 steps and compared it to PGD on
25 the same subset. Since these attacks were extremely computationally expensive, we were only able to run them on 10
26 images for the 10k steps. So far, the trend agrees with the observation made in the paper (see Figure below). In all
27 cases, the size of the perturbation was smaller with PGD, while the probabilities assigned to the mispredicted class for
28 examples generated with 20 steps PGD were much greater than for 10k steps of boundary attack.



29 **(R3) Marginal improvement.** The reviewer critiques that our proposed mechanisms only improve the robustness
30 marginally as compared to standard training. While this is true for large perturbations (which both R3 and we note
31 suggests additional mechanisms), the effect at small perturbations is strong: at $\epsilon = 0.005$ the adversarial performance
32 more than doubles over the baselines both in CIFAR10 and ImageNet10.

33 **(R1) Comparison to adversarial training.** While the improvement from both mechanisms is lower than that of
34 adversarial training, at $\epsilon = 0.005$ on CIFAR10 they reach approximately the halfway point between standard models
35 and those adversarially trained, as seen in Figure 3 in the paper. However, unlike adversarial training, which is very
36 computationally expensive, our methods pose negligible computational overhead and actually lead to improvement of
37 standard performance on some of the datasets with the retinal fixations – this should be compared to adversarial training,
38 which provably leads to significantly worse standard performance. Additionally, we have shown in the Supp. Mat. that
39 our methods can scale to larger datasets (see the ImageNet100 experiments), while so far adversarial training has been
40 unsuccessful in being scaled up to larger number of classes (hence the suggestion of R2 to try adversarial training on
41 ImageNet is beyond the scope of this work).