

1 We thank the reviewers for their constructive comments. We are pleased that even the most critical reviewers found  
 2 the claims to be of "substantial significance" if they are validated with key controls. *We have completed all of those*  
 3 *controls*; while some were in the original submission but not sufficiently highlighted, others are new analyses following  
 4 reviewer suggestions. In all cases, the controls support our claims and we would happily include these improvements in  
 5 the final manuscript. We thank the reviewers for their willingness to reconsider their reviewer score in this light.

6 **VOneBlock describes V1 responses better than CNNs.** We agree with R1 and R3 that  
 7 if the VOneBlock is not better than CNNs at describing V1 responses, it would undermine  
 8 our claims. In lines 155-158 we write "the VOneBlock outperformed all tested ImageNet-  
 9 trained CNNs in explaining responses in the V1 dataset used", and provided the value  
 10 of explained variance ( $0.387 \pm 0.007$ ). Here, we update Fig.1, showing VOneResNet50  
 11 surpasses all standard ImageNet-trained CNNs. We are confident in this result since we  
 12 searched over all layers in a large pool of CNNs and did not find any coming close to  
 13 the V1 explained variance of the VOneBlock. Further, **our results are consistent with**  
 14 **Cadena et al 2019**: our GFB has parameters constrained by empirical data resulting in  
 15 a better model of V1; when we use the parameters of the GFB in Cadena et al, we obtain  
 16 a much lower explained variance ( $0.296 \pm 0.005$ , marked on x-axis of Fig.R1).

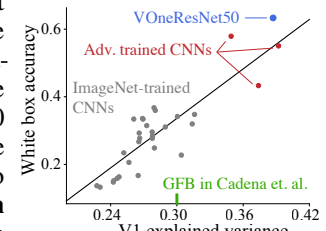


Fig R1: Accuracy in weak white box attack vs V1 explained variance; VOneBlock outperforms standard CNNs

17 **Stochasticity during the attack is not the main source of adversarial robustness.** R1 and R4 are concerned that  
 18 stochasticity makes VOneNets artificially appear more robust. However, the majority of the VOneNets' robustness does  
 19 not originate from stochasticity during the attacks, but rather from training downstream layers with V1-like features and  
 20 neural stochasticity. *When we turn off stochasticity completely during attack/inference, VOneResNet50 retains 70% of*  
 21 *the adversarial robustness gains of the model with stochasticity on, and significantly outperforms ResNet50* (Fig.4).  
 22 Meanwhile, when we add stochasticity to ResNet50, using an affine transformation to scale the activations to the same  
 23 magnitude as in the VOneBlock (lines 148-151), we find substantially smaller improvements in robustness (Fig.3B).

24 **Further adversarial attack optimization does not overturn results.** We follow  
 25 R1's suggestions to further verify our attacks, addressing R4's concerns of gradient  
 26 quality as well. Increasing the attack step size ( $\alpha$ ) and the number of gradient  
 27 samples ( $k$ ) only marginally increases the attack effectiveness. A grid search over  
 28  $\alpha$  and  $k$  (Fig.R.2; PGD with  $\|\delta\|_\infty = 1/255$ ) reveals that at the most effective step  
 29 size ( $\alpha = \epsilon/8$ ), increasing gradient samples beyond  $k = 16$  no longer improves  
 30 the attack. At  $k = 128$ , VOneResnet50 accuracy is only reduced from 29.1% to  
 31 26.0%, remaining a large margin above ResNet50 at 0.8%. We do not find using  
 32 5-10 random restarts improves beyond the most effective  $k$  and  $\alpha$  settings. Thus,  
 33 these controls do not qualitatively change our results. Finally, because our model  
 34 has a stochastic boundary, we focus on PGD and avoid boundary targeted attacks,  
 35 as the most effect should come from staying as far from the boundary as possible.

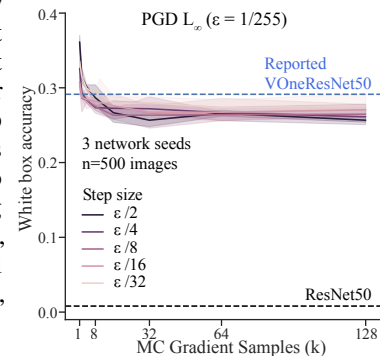


Fig R2: Attack optimization only marginally decreases VOneResNet50 accuracy

36 **Common corruptions provide additional insights.** R1 and R3 question the in-  
 37 clusion of common corruptions. While common corruption scores generally do not  
 38 correlate with V1 similarity, we include them because robustness to imperceptible adversarial attacks should not come  
 39 at the cost of performance on other perturbations that humans easily deal with, as is the case for adversarially trained  
 40 models which underperform on this benchmark relative to the baseline and VOneNets. Further, common corruptions  
 41 provide insights into the role of different components in the VOneBlock in dealing with specific stationary perturbation  
 42 statistics. **Common corruptions and white box attacks are properly separated:** results are shown separately in  
 43 Tables 1, 2, B.2, C.4, and C.5; white box only in Fig.1, 2, and 4; and scores are only combined in Fig.3.

44 **Mechanisms contributing to robustness.** We agree with R1 and R3 that the mechanisms behind our results warrant  
 45 further exploration, but believe our study provides a novel and important characterization of how low level visual  
 46 processing relates to robust object recognition behavior. With ablations, we characterize how different components of  
 47 the VOneBlock impact performance on a variety of perturbation statistics. Lines 212-218 discuss the effect of high  
 48 SF Gabors on various corruptions as well as complex cells on adversarial attacks, and Fig.4 shows that the improved  
 49 adversarial robustness largely stems from training downstream layers with a stochastic V1 front-end.

50 **Other questions from R1 and R3.**  $\pm$  refers to SD. As stated in the appendix, we will release code and model weights.  
 51 To the best of our knowledge, the correlation between adversarial robustness and V1 similarity has not been shown  
 52 before. The CNN layer that best predicts V1 responses is usually not the first (except for the VOneNets). We used  
 53 parameters from the foveal region of V1 (eccentricity  $< 5$ deg). Our model has a field-of-view of 8deg (resolution =  
 54 224px/8deg; eccentricity  $< 4$ deg). The VOneBlock and all tested CNNs explained less than 40% of V1 response variance  
 55 which is in line with the literature and is a possible leeway to make progress on robustness. One of the comparison  
 56 models (ANT3 $\times$ 3+SIN) is trained with input noise and does not show improvements in adversarial robustness.