

1 Thank you for yours detailed reviews. We would first like to address the **general concerns**. Firstly, we agree that the  
 2 clarity and writing of the paper needs improvement, both in terms of notation and explanation of various terms, and  
 3 especially in the section on adversarial attack detection. We have taken your comments on board and are currently  
 4 addressing these issues in order to make the story of the paper far clearer and easier to follow. Secondly, we agree that  
 5 the experimental evaluation needs improvement. This was mostly due to having legacy tensorflow code that was not  
 6 suitable for extension. The last few months have been spent writing a cleaner PyTorch implementation which allows for  
 7 the easier introduction of new architectures and integration with adversarial toolkits (FoolBox/Cleverhans) and allows for  
 8 a head-to-head comparison of various methods. We are currently reproducing (classification accuracy results presented  
 9 in table below) the OOD detection experiments using the WideResnet 28x10 architecture, where we are able to match  
 10 SOTA classification performance on CIFAR-10, CIFAR-100 and TinyImageNet. Additionally, we will add experiments  
 11 on RKL PNs trained on TinyImageNet in-domain and a 400-class subset of the other 800 imageNet classes, processed  
 12 like TinyImageNet, as OOD training data. We will evaluate OOD detection on the remaining, heldout subset of 400  
 13 ImageNet classes. Third, we will update adversarial attack detection numbers to be on the CIFAR-10, CIFAR-100 and  
 14 TinyImageNet datasets using the new architecture. Furthermore, having analyzed the Carlini&Wagner L2 attack, we  
 15 believe that the current adaptive adversarial attack loss functions may be sufficient. An alternative adaptive attack loss  
 16 function we will also consider is L1 loss between the predicted and permuted logits, which has the same fixed point as  
 17 KL-divergence minimization, but potentially nicer gradients. This should yield an evaluation against a stronger adaptive  
 18 whitebox adversary. However, we are unable to demonstrate updated Adversarial Attack Detection numbers in this  
 rebuttal as we haven't integrated with Foolbox yet.

Model	CIFAR-10	CIFAR-100	TinyImageNet
VGG (Error %)	8.0 $\pm$ 0.4	30.4 $\pm$ 0.6	41.7 $\pm$ 0.4
WRN (Error %)	3.9 $\pm$ 0.1	19.3 $\pm$ 0.5	32.3 $\pm$ NA

19

20 **Reviewer 2** The OOD results quoted in the papers use a range of classification network architectures. It is thus not  
 21 appropriate to directly compare numbers. We are planning to do clean comparisons of the method using consistent  
 22 architecture. Secondly, as presented in this paper, Prior Networks are *not* SOTA (but close) for OOD detection, as  
 23 the current SOTA (Lee18) results make use of a set of *bespoke* post-processing techniques (such as ODIN) aimed at  
 24 maximizing OOD detection performance. The same approaches can be applied on top of Prior Networks. The aim of this  
 25 paper is to present a *general method* (training criterion) and analyze its properties in two different scenarios. Crucially,  
 26 we show that adversarial training of Prior Networks using the proposed criterion gives a significant improvement over  
 27 standard adversarial training *at not additional cost*, making it a drop-in replacement. Other techniques can be stacked  
 28 on top of this. Third, a large perturbation size was selected in order to give more freedom to the adversarial attack  
 29 to succeed against our detection scheme. In additional results to be completed we will conduct an analysis of the  
 30 perturbation size on the success of adaptive whitebox attacks. Fourth, perhaps we have misunderstand your comment,  
 31 but the joint success rate represents the success of *the attack*, rather than *the defense*. Specifically, the JSR represents  
 32 the success of the attack at *both* successfully attaining the target class *and* avoiding detection. This will be made clearer  
 33 in the text. Thus, it is unsurprising that an adaptive whitebox adversarial attack will have a very high success rate.  
 34 Fifth, regarding the dropout attack. We need to make it clear in the text that attacks against dropout are completely  
 35 undetectable, but have a slightly lower success rate since we generate the attack against the *mean network*, rather than  
 36 against each of the 10 samples. The added stochasticity makes the attack less successful. This was difficult to address  
 37 in the legacy code, but should be easier to do in the new implementation. Sixth, while we agree that the computational  
 38 expense of an adversarial attack can be increased by gradient masking, we think that it is a significant result that we can  
 39 do the same by essentially using an improved form of adversarial training at no additional expense. We stress that we  
 40 analyse a general method. Other, more task-specific techniques can always be stacked on top of this. Finally, we use  
 41 FGSM attack in training because we dynamically generate the attack on each minibatch during training. Using iterative  
 42 attacks is possible, but expensive, as training would slow down.

43 **Reviewer 3** Firstly, the Reverse KL is exactly equal to the variational criterion, where the reconstruction loss is weighted  
 44 by  $\beta$  and where the KL-regularization loss is minimization the KL-divergence to a flat Dirichlet Prior. We will make  
 45 this connection clearer in the paper. Secondly, this choice of OOD training data is quite standard in OOD detection  
 46 experiments. The main requirement is that it is more diverse than the in-domain data. Third, the effect of varying  $\beta$  is  
 47 not strong. The main goal is to make sure that the distribution of  $\hat{\alpha}_0$  in-domain and OOD are clearly separable. It is  
 48 necessary to set  $\beta$  to 0 for OOD training data (for OOD detection), as we have no a-priori knowledge about what the  
 49 target class. However  $\beta$  can be set to 1 for adversarial data, as we do know what the target class should be. This allows  
 50 the model to learn to *both* predict the correct target class *and* high uncertainty for adversarially perturbed training data.  
 51 Finally, setting  $\gamma = 10.0$  for the RKL loss was chosen based on results of the toy-data experiments. When training  
 52 models on the other datasets we found that this choice of gamma was consistently better than setting gamma to 1.0 .