

1 **Reviewer 1** : Thanks for your positive and constructive feedback. We address your detailed comments below. **Re:**  
2 **difference between logits yielding lower adversarial accuracy than the CE:** Optimizing the difference between  
3 logits is very similar to the Carlini-Wagner attack [1]. We tried optimizing both cross entropy and the difference  
4 between logits, and found the latter to be a stronger attack. **How is the regularization parameter chosen:** Thanks,  
5 these details were indeed missing. We do a parameter sweep; we will provide full details in the appendix. **Can**  
6 **TRADES outperform the proposed method:** The performance of TRADES reported in the paper was obtained by  
7 taking the max of results from an extensive sweep over weights on the regularizer. We will clarify and incl. full sweep  
8 results in the appendix. **SOTA results on ImageNet are easier to achieve than for CIFAR-10, MNIST / Compare**  
9 **to more methods:** Note that we match the SOTA on CIFAR-10. The main point of the paper, however, was to develop a  
10 method that would **scale to ImageNet** (which previous methods found difficult). We put significant effort and time into  
11 carefully creating two strong baselines (Adversarial Training and Denoise), tuning them extensively for ImageNet and  
12 comparing all methods under the same attack and using the same network architecture. **Most groups do not have 128**  
13 **TPUs:** We would like to note the efficacy of our method is not due to the amount of available compute. It outperforms  
14 competing approaches even in the low-compute regime. It can also be run on GPUs and it would only be 2x more  
15 expensive than standard ImageNet training. In comparison, adversarial training is 30x more computationally expensive.

16 **Reviewer 2:** We thank the reviewer for the feedback and address the raised questions below. We hope that these answers  
17 clarify points that were unclear and will revise the paper accordingly. **No work verifies that preventing gradient**  
18 **obfuscation leads to better robustness:** Our paper shows that by maintaining locally linearity we enforce robustness  
19 (see Section 4.3 and Appendices C, D). Moreover, empirically we have observed local linearity also avoids gradient  
20 obfuscation when we train with much fewer steps of PGD than adversarial training (see next question) and we hence  
21 here make a connection between the two. As the reviewer mentioned, no work has made this connection before. We  
22 will try to make this point clearer. **Why is minimizing the upper bound better than directly minimizing the loss**  
23 **gap (PGD-training):** We agree with the reviewer that if we could consistently find the optimal attack that maximizes  
24 the loss gap using PGD then training with this attack would be effective. However, to find a sufficiently strong attack  
25 for large models would require a significant amount of compute (e.g. 30 steps of PGD on ImageNet), while training  
26 with fewer PGD steps can lead to gradient obfuscation (see Section 3.2 and 5.3). The motivation for using LLR is to  
27 encourage a linear loss surface and thus prevent gradient obfuscation. By enforcing local linearity, our regularizer makes  
28 it easier to find a strong attack with much smaller number of PGD steps. Indeed, if the loss surface is linear, then PGD  
29 can find the optimal attack in a single step. **Re: results on more diverse attacks, like Deepfool:** Rather than a diverse  
30 set of attacks we found it more important to choose the strongest attack and compare different baselines in the same  
31 framework under consistent attacks. As demonstrated by Carlini and Wagner [1] (and confirmed by TRADES [2]) their  
32 attack is much stronger than DeepFool. Note: we also devised a stronger attack (Multi-Targeted Attack) which achieves  
33 the lowest adversarial accuracy. We do, however, understand the concern and will include DeepFool. Preliminary results  
34 (for WRN-28) TRADES: 63.49%, LLR: 71.43%. **Regarding more convincing motivation and more experiments:**  
35 We hope the above addresses your concerns regarding the motivation. We believe we presented an extensive set of  
36 experiments on CIFAR-10 and ImageNet (re-implementing the baselines for equal comparison with strong attacks).  
37 We further investigated the change in accuracy as we increase the strength of attack for both LLR and all baselines.  
38 Moreover, we ablated different parts of the regularizer. We also performed a statistical analysis on the linearity measure  
39 comparing different adversaries to our LLR. Much of this experimentation is in the appendix, but they are referred to in  
40 the main text. If there is an experiment missing we are happy to include it.

41 **Reviewer 3:** We thank the reviewer for the feedback, especially regarding mathematical details. It is appreciated. We  
42 address the comments below. **Re: (Eq. (8)). As we seek to minimize  $\langle \delta, \nabla_x \ell(x) \rangle$  for all perturbations  $\delta$  in the local**  
43 **neighborhood  $B_\epsilon$ , we should naturally aim at minimizing  $\|\nabla_x \ell(x)\|_2$ . Why use  $\langle \delta_{LLR}, \nabla_x \ell(x) \rangle$  instead ?:** It's  
44 true that  $\langle \delta, \nabla_x \ell(x) \rangle \leq c \|\nabla_x \ell(x)\|_2$  for  $c = \max_{\delta \in B_\epsilon} \|\delta\|_2$ , and thus if we wanted to minimize  $\langle \delta, \nabla_x \ell(x) \rangle$  for all  
45  $\delta$  then  $\|\nabla_x \ell(x)\|_2$  is a good objective to minimize. We have tried this but found this bound to be less effective in  
46 practice. Concretely, if the weight on  $\|\nabla_x \ell(x)\|_2$  is small then it is not better than training with  $\gamma(\epsilon, x)$  alone (49.37%  
47 adversarial accuracy), if large it has significant impact on the nominal accuracy (reduction to 80%). This could be due  
48 to the fact that  $\|\nabla_x \ell(x)\|_2$  is a looser bound than the one we optimize; and constrains the rate of change of the loss  
49 in all directions. We are happy to include these observations in an updated version of the paper. **Compare running**  
50 **time of the proposed method to that of CURE:** Thanks for this comment; we should have mentioned this and will  
51 clarify in the paper. The running times' comparison is as follows: CURE is essentially performing 2-steps of GD (to  
52 approximate the curvature); thus our ImageNet running time is the same. For CIFAR-10 we optimized for robustness  
53 and not training time. We could have done 2-step PGD – see appendix – but 15-steps gives better results.

## 54 References

- 55 [1] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *IEEE Security and Privacy*, 2017.
- 56 [2] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled  
57 trade-off between robustness and accuracy. *arXiv preprint arXiv:1901.08573*, 2019.