

1 We thank reviewers for their thoughtful feedback. We are pleased to see that most reviewers found the work to be  
 2 interesting and innovative. Here, we provide clarifications and conduct additional experiments to demonstrate the  
 3 improvement in performance achieved by our proposed technique for certification under the  $\ell_\infty$  norm threat model.

4 **Reviewer 1 : 1) Significance of results.** Although we only give empirical evidence for  $\ell_1, \ell_2$  norm and subspace  $\ell_2$   
 5 norm in the manuscript, the theoretical guarantees in fact extend to  $\ell_\infty$  and subspace  $\ell_1, \ell_\infty$  norms. We only focus on  
 6  $\ell_1, \ell_2$  and  $\ell_\infty$  as these are the most intensely researched / relevant threat models in the field. In response to Reviewer  
 7 3’s comments, we also provide empirical evidence to show improvement for  $\ell_\infty$  norm on the CIFAR10 dataset. As  
 8 the current state-of-the-art for  $\ell_\infty$  norm is given under Gaussian smoothing [1], our empirical result can give the new  
 9 state-of-the-art for  $\ell_\infty$  norm certification. **2) Cohen  $\ell_1$  radius calculation :** The  $\ell_1, \ell_\infty$  norm results were not explicitly  
 10 stated in the original paper [2] but they can be derived by following the same analysis, which are also stated in the recent  
 11 paper [1, Appendix Table A] . **3) "Higher-order" in title :** The "higher-order" in the title is in reference to the fact that  
 12 the paper lays down the ground work for using higher-order information for certification. However, we see that it might  
 13 be ambiguous as we only fully explore first-order smoothing. So, we plan to change the title to "first-order smoothing".

14 **Reviewer 2 : Limited empirical evidence :** We note that our experiments provide numerical evidence of our theoretical  
 15 results and demonstrate that the certification performance can be greatly improved by incorporating higher-order  
 16 information. We have followed standard experiment setup and conducted various experiments on CIFAR10 (Sec 5) and  
 17 ImageNet (Appendix E) and compared all the current baselines the  $\ell_1, \ell_2$  norm and subspace  $\ell_2$  norm, which is in line  
 18 with other works in this field. Notably, we conduct additional CIFAR10 experiments for  $\ell_\infty$  certification in Figure R1.

19 **Reviewer 3 : Experiments for  $\ell_\infty$  :** In the current paper, we  
 20 have only focused on giving the bounds for  $\ell_1, \ell_2, \ell_\infty$  norms.  
 21 For general  $\ell_p$  norm we can use the current results to provide  
 22 lower bounds on the certified radii. As for empirical results  
 23 for certifying the  $\ell_\infty$  norm radius, it requires the estimation  
 24 of  $\|y^{(1)}\|_1$ . As mentioned in line 270 in the paper the current  
 25 estimators used to calculate  $\|y^{(1)}\|_1$  need a lot of samples  
 26 in order to find non-vacuous high-confidence bounds. Although  
 27 the certification cost is higher, **using the proposed method**  
 28 **gives us significant (~ 10%) improvement over the bounds**  
 29 **given by Cohen et al. for CIFAR10.** The CIFAR10  $\ell_\infty$  results  
 30 in Figure R1 are calculated using 4M samples for certification  
 31 (6 minutes/image). One of the major limitations of the current  
 32 estimators is that they are biased. In the paper we have proposed  
 33 a new unbiased estimator for  $\|y^{(1)}\|_2$  (Table 1 in the current  
 34 paper). We think a similar new estimator is needed to make  $\ell_\infty$   
 35 certification more scalable. This is left for future work.

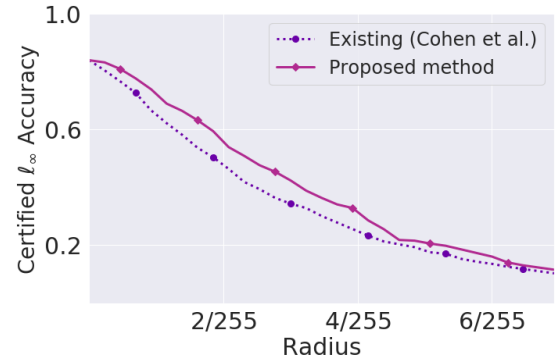


Figure R1: Comparing certified accuracy for CIFAR10 under  $\ell_\infty$  threat models. Our results show that around 10% improvement can be obtained by using the proposed method.

36 **Reviewer 4 : 1) Typos and clarifications :** **i)** Sorry it is a major typo. The statement should read "*Under the proposed*  
 37 *general framework for calculating certified radii, it is easy to see that adding more local constraints ( $H_i^x$ ) in Equation*  
 38 *(2) gives a smaller bigger value of  $p_x(z)$  for any  $x, z$  which makes the super-level set of  $p_x$ , equivalently the certified*  
 39 *safety region, bigger."* **ii)** With a slight abuse of notation, we use  $\mu$  to denote both the measure and the probability  
 40 density function of the measure. Here,  $D_x^\alpha \mu(y-x)$  corresponds to taking the multivariate differential of the probability  
 41 density function ( $\mu$ ) at  $(y-x)$  with respect to variable  $x$ . **iii)** In Figure 1 of our manuscript, the direction of the gradient  
 42  $y^{(1)}$  is along the negative x-axis. We plan to add an arrow to clarify this. Also in order to better motivate the idea we  
 43 plan to add numerical values for the directional certified radii on the figure. **iv)** Given the images in the pixel space, we  
 44 do change of basis to orient the basis along the gradient  $y^{(1)}$  to simplify calculations. In line 515-518,  $z_1, z_2, \dots, z_d$   
 45 denotes the variables corresponding to the new basis vectors for the space after the transformation. In corollary 1 we  
 46 abuse the notation and use  $z_1, z_2$  to denote the variables involved in the system of equations we reduce our initial  
 47 constraints to. The two sets of variables are not linked. We will change the variable names to avoid confusion in the  
 48 future and also give a description of  $z_i$ 's in the proof before using them. **2) Certified bounds vs Attack bounds :** We  
 49 do agree that the experimental evidence would be great. However we are aware of attacks on randomized smoothing  
 50 classifiers only for the  $\ell_2$  norm threat model currently. For this scenario the current bounds are near optimal as the  
 51 attacks are close to the current state-of-the-art certification bounds [3] (equivalently our proposed certification bounds).

52 **References**

54 [1] G. Yang, T. Duan, E. Hu, H. Salman, I. Razenshteyn, and J. Li, "Randomized smoothing of all shapes and sizes," *ICML*, 2020.  
 55 [2] J. M. Cohen, E. Rosenfeld, and J. Z. Kolter, "Certified adversarial robustness via randomized smoothing," *ICML*, 2019.  
 56 [3] H. Salman, J. Li, I. Razenshteyn, P. Zhang, H. Zhang, S. Bubeck, and G. Yang, "Provably robust deep learning via adversarially  
 57 trained smoothed classifiers," *Neural Information Processing Systems*, 2019.