

1 We thank the reviewers for their insightful and constructive feedback. We start with addressing two common issues  
2 (**C#1**, **C#2**), and then the comments from the individual reviewers (**R1**, **R2**, **R3**, **R4**).

3 **C#1: About the connection with [A]**. We will definitely draw more comprehensive discussions on [A] as suggested  
4 by **R3**. We believe the two works have different contributions and were developed concurrently. [A] derives the  
5 optimal shapes of level sets for  $\ell_p$  attacks based on the Wulff Crystal theory, while our work, based on our functional-  
6 optimization framework and accuracy-robustness decomposition (Eq.9), proposes to use distribution that is more  
7 concentrated toward the center. Besides, we also consider a novel distribution using mixed  $\ell_2$  and  $\ell_\infty$  norm for  $\ell_\infty$   
8 adversary, which hasn't been studied before and improve the empirical results.

9 **C#2: About the significance of empirical results**. As we simply use the pre-trained model from open source of the  
10 baselines and thus our algorithm is almost deterministic (as the variance of Monte Carlo is very small given that we use  
11 100K samples). Our result are statistically significant, for example for the 49% vs. 50% accuracy for ImageNet on  $\ell_2$   
12 adversarial, two-proportions z-test gives p value 0.001565.

13 **R1 On  $\ell_\infty$  results**. We achieve better or comparable performance as IBP. We will add the comparison in revision.  
14 Previously, we do not list IBP results because it relies on a special model structure while our method can be applied to  
15 any model structure. [31] uses trades to train better model, while we focus on better certification accuracy with a given  
16 model. Applying our method to [31]'s checkpoint (image size = 32,  $\epsilon = 0.435$ ,  $\beta=6$ , standard deviation = 0.12), we  
17 achieve 64.7% accuracy at  $2/255$  radius outperforming 62.6% in [31]. We will add this baseline in the next version.

18 **R1 On the optimality of Gaussian**. Yes, [30-32] gives the (rate) optimality of Gaussian for  $\ell_p$  adversarial. But this  
19 theoretical results only analyzes the rate of the error while ignoring the constant. However, in practice, the constant can  
20 be important or even dominating (e.g. many non-convex optimization algorithms have same rate while can be quite  
21 different in practice due to the difference in constant factor).

22 **R1 On the connection between the dual form and proposed distribution**. The proposed distribution family is the  
23 one that gives empirical improvement. The theoretical framework serves as the motivation (decomposition of the bound  
24 in sec 3.3) for proposing the new distribution family. It is not valid to do ablation study on the theoretical framework  
25 (i.e. when we choose previous smoothing distribution, i.e., Gaussian in [C], our theoretical framework gives same result,  
26 see Corollary 1,2).

27 **R2 [B] and our work are two different ways for generalizing the original random smoothing bounds**: [B] relaxes the  
28 constrain of smoothing samples to a ball in f-divergence while ours relaxes the classifier to be in a function space. We  
29 will draw an in-depth discussion in the revision. Regarding the choice of  $k$ , we refer readers to Appendix B.3 for details.  
30 We will move the reject sampling details to main text.

31 **R3 On the 'weakness'**. Thanks for the information on [A]. Regarding work [A], please see C#1. For the questions  
32 on  $\ell_2$  adversarial, our improvement does not come from the optimization as we simply use the pretrain model from  
33 open source repository of the baselines. We suspect the empirical different comes from the different choice of hyper-  
34 parameters of the distribution. We are not sure what you mean by saying 'analytic form gives different findings'. Could  
35 you explain it with more detail so that we can address it in the next version? For the question on  $\ell_1$  adversarial, yes, [A]  
36 achieves better empirical results. But our results support our theoretical findings and is of independent values to the  
37 community. We will definitely add comprehensive discussion on [A] in the next version.

38 **R3 On the minor remark about 'correctness'**. Thanks for pointing this out. We agree using sup/inf is more precise  
39 and will change to it. And yes,  $\mathcal{F}_{[0,1]}$  isn't compact. We will fix this. We will also add proof for theorem 1 (III).

40 **R4 On the trade-off**. In line 143-147, we explain our main insights gained from Eq.9, from where we carefully design  
41 the proposed distribution to achieve a better trade-off (see Fig.1). The  $\lambda$  comes from the dual derivation, and its optimal  
42 value only depends on the specific value of the accuracy term and the norm ball in Eq.9. Even though the optimal  
43  $\lambda$  is generally hard to derive analytically, the two components in Eq.9 clearly show the trade-off between accuracy  
44 and robustness. Similarly, performing rigorous scaling analysis on  $\ell_\infty$  is rather challenging technically. We leave this  
45 analysis for future work.

46 **R4 On empirical result**. Please see C#2 on the statistical significance. The tightness of our bound is proved in theorem  
47 1 (III). Empirically, the exact certified accuracy relies on  $f_{\pi_0}^\#$ , which depends on whether a good smoothing distribution  
48 that fits the data and base classifier are selected (which is also the main gain of this paper).

49 **R4 On adversarial training**. Thanks for pointing out this. We agree it's doable while it would require an unaffordable  
50 amount of computational resource for adv. training on ImageNet. Thus we leave this for future work.

51 **Reference [A]** Randomized Smoothing of All Shapes and Sizes. **[B]** A Framework for Robustness Certification of  
52 Smoothed Classifiers Using f-Divergences. **[C]** Certified Adversarial Robustness via Randomized Smoothing.