

1 We thank the reviewers for their comments. We address individual concerns below. If you think we address your
2 concerns, please consider raising the score.

3 **Reviewer 1: Simple method and limited contribution.** We believe the simplicity of our method is a strength, not a
4 weakness. In addition, it's not enough to just "plug and play" Cohen et al.'s method, as our No-Denoiser baseline shows
5 (Tables 1&2 and Figure 2). Thus our algorithmic contributions, though simple, are important.

6 *Relationship with prior work.* All prior work which apply a preprocessing step are empirical defenses that hope to
7 remove malicious perturbations by doing a preprocessing step. In our work, we apply a denoiser not to remove the
8 malicious noise, but to make the pre-trained classifier accurate under Gaussian perturbation of its input, therefore
9 making randomized smoothing effective when applied to this pre-trained classifier.

10 *The white-box scenario is misleading.* We believe Reviewer 1 misunderstood the meaning and context of "white-box"
11 in our paper. Whether the attacker has access to the denoiser or not doesn't affect our method as we aren't empirically
12 removing adversarial noise. Our guarantees are provable due to the use of randomized smoothing.

13 *Practicality: Performance gap w.r.t to Cohen et al's method.* One should expect our method, without retraining the
14 classifier, performs at most as well as Cohen et al. which trains the classifiers with smoothing in the loop. While the
15 gap between these methods point to important future work direction, the value of our contribution is clear by comparing
16 against the No-Denoiser baseline, which is the *real "plug-and-play"* referred to by the reviewer in point 1. See Tables
17 1&2 and Figures 2,3,&4.

18 *What threat models can the algorithm handle?* Since our method uses randomized smoothing, it can in theory handle all
19 threat models that randomized smoothing can handle (including l_1 , l_2 , l_∞ , and Wasserstein). The only change in our
20 method would be the way our denoisers are trained. For instance, instead of training a denoiser to remove Gaussian
21 noise (for L_2 certification), the denoiser shall be trained to remove noise sampled from other distributions (e.g. Laplace
22 distribution for L_1 threat models). We agree this might be confusing in our paper so we will make sure to clarify it.

23 *No proper theoretical justification.* We are not sure what theoretical justification the reviewer is pointing to here. Our
24 method applies randomized smoothing to a composition of a classifier and denoiser, instead of only applying it to
25 a classifier as all prior works on randomized smoothing do. Therefore, all the theoretical guarantees of randomized
26 smoothing hold for us.

27 *The paper is tough to follow and read.* We will reorganize the paper to be more self-contained in the main text.

28 *Reproducibility* We provide detailed experimental details in Appendix B, along with a detailed code + pre-trained
29 models replicating all the experiments. So we are confused why the reviewer thinks the paper is not reproducible.

30 **Reviewer 2: The authors break the promise of avoiding any re-training that is given in the paper.** We stress that we
31 never re-train the base classifier neither in the white-box nor black-box settings. In the white-box setting, we assume
32 we know the base classifier, and we backpropagate through it, but we only update the denoiser. The whole purpose
33 of our paper is to get non-trivial certification results without re-training the classifier itself. We hope this clarifies the
34 confusion; we will update the paper accordingly.

35 *Comparison to Madry's adversarial training.* Madry's defense is empirical, whereas we are interested in certified
36 defenses. But in any case, Madry's defense requires adversarially training the classifier, whereas in our setting, the
37 classifier isn't allowed to be re-trained/modified at all. It is interesting to study whether PGD-like defenses can be
38 applied to pre-trained classifiers without modifying the latter, but this is outside the scope of this paper.

39 *Limited novelty.* We agree that our approach looks similar to the use of denoising auto-encoder from Lecuyer et al.
40 However, our approach is distinguished from Lecuyer et al in several ways:

- 41 1. The use of denoising auto-encoder in Lecuyer et al is solely aimed at speeding up training for certifying large models.
42 In our case, our motivation is to effectively apply randomized smoothing to pretrained models.
- 43 2. More importantly, we comprehensively investigate various training strategies (MSE/stability/classification objectives)
44 and application settings (white-box/black-box), which are not investigated in Lecuyer et al.

45 *Practicality issues.* Our approach utilizes randomized smoothing, thus any practicality issue of randomized smoothing
46 passes on to our approach. We agree with the reviewer on this point and we will make it clearer in the next revision.

47 *Gap between ours and Cohen et al. increases as the model becomes more complex.* This is indeed an interesting
48 observation that needs further investigation. It might be the case that larger architectures require different denoiser
49 architectures/training schemes. We don't claim we train our denoisers in the best way possible, and we believe with
50 improved training of the denoiser, the gap can be further reduced.

51 *Reconstruction artifacts of STAB+MSE compared to MSE* We think the reason for this is that STAB+MSE makes the
52 denoiser more customized to the base classifier, resulting in more corrupted reconstruction. MSE loss only considers
53 removing Gaussian noise, thus leading to visually better output. This requires further investigation and is left for future
54 work.

55 **Reviewer 3:** Thanks for your comments!