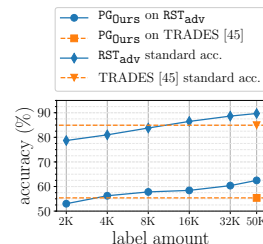1  We thank the reviewers for the kind and helpful reviews. Below we address each review in turn.

2  **Reviewer 1**  We thank the reviewer for recognizing the novelty and significance of our
3  results, as well as for the additional important references about robust generalization—
4  we will include them in the revised paper. We also thank the reviewer for the excellent
5  suggestion to study how many *labels* are actually necessary for achieving state-of-the-art
6  robustness. Following this suggestion, we perform robust self-training with random subsets
7  of CIFAR-10 as the labeled data and the remainder of CIFAR-10 and our mined 500K
8  images from Tiny Images as unlabeled data. The figure to the right shows the result (for
9  adversarial training and testing with $\mathrm{PG}_{\mathrm{Ours}}$ with $\epsilon = 8/255$): with as few as 4K labels, we
10  are able to match the fully-supervised state-of-the-art! The revised paper will include a
11  detailed account of this experiment.



12  **Reviewer 2**  We thank the reviewer for the valuable feedback, which will improve the readability of our paper. Below,
13  we address each point in the review; we will also carefully revise our paper to clarify each of these points.

14  *"The term $\ell_2$ certified accuracy is not defined."* By certified $\ell_2$ accuracy $\xi$, we mean a proof that $\mathrm{err}_{\mathrm{robust}}^{2,\epsilon}$ defined in Eq. (2)
15  is at most $1 - \xi$ on the test set. Specifically, we use the randomized smoothing proof by Cohen et al. (2019).

16  *"Line 94 says difficult to learn classifier. For the Gaussian model, the classifier must be easy to learn. Isn't that so?"*
17  Here by "difficult to learn" we meant "requires many samples to learn."

18  *"Line 126-128. It is difficult to follow the logic..."* Here is a more detailed explanation: we have $\hat{\theta}_{\mathrm{final}} = (\frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \tilde{y}_i y_i)\mu +$
19  $\frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \tilde{y}_i \varepsilon_i$ where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2 I)$ is the noise in example $i$. In the proof we show that with high probability
20  $\frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \tilde{y}_i y_i \geq \frac{1}{6}$ while the variance of $\frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} \tilde{y}_i \varepsilon_i$ goes to zero as $\tilde{n}$ grows, and therefore the angle between $\hat{\theta}_{\mathrm{final}}$ and
21  $\mu$ goes to zero. By Eq. (11) in Appendix A.1 and Eq. (3) this implies that the robust error becomes small.

22  *"Using more datasets for experiments might be more convincing."* We agree and will gladly experiment on more datasets.
23  Unfortunately, there are not many established benchmarks for adversarial robustness (we do not have the computational
24  resources for adversarial training on ImageNet). We would appreciate suggestions for additional datasets to consider.

25  *"What happens when there is class imbalance?"* The theoretical results in this work easily extend to the case where there
26  is class imbalance. The upper bounds in Proposition 1 and Theorem 2 hold regardless of the label distribution. The
27  lower bound in Theorem 1 changes from $\frac{1}{2}(1 - d^{-1})$ to $p(1 - d^{-1})$ where $p$ is the proportion of the smaller class; the
28  only change to the proof in [34] is a modification of the lower bound on $\Psi$ in page 29. We thank the reviewer for raising
29  this interesting question—we will mention it in the revised paper.

30  **Reviewer 3**  Thanks for the helpful comments (addressed below) and for finding our paper a pleasure to read.

31  *"The main concern is that the connection between the theory and the experiment is loose."* We believe that there is a
32  substantial connection between theory and experiment in our paper because both parts follow the same algorithmic
33  approach. As we mention in lines 288–289, this is not a coincidence: our theoretical results motivated our experimental
34  investigation. In particular, the observation that self-training is very effective in utilizing unlabeled data in the Gaussian
35  model led us to empirically test it in more realistic settings. We agree that showing a sample complexity separation in a
36  more realistic model is a challenging open problem. We thank the reviewer for pointing out that this connection wasn't
37  clear enough in our paper; we will state it clearly in the introduction to the revised paper.

38  *"The comparison seems to be unfair with the state of art models because robust self-training has extra unlabeled data*
39  *information."* Our main contribution is to show that unlabeled data improves adversarial robustness, and therefore our
40  primary experiments focus on evaluating this improvement using state-of-the-art models. However, in Table 1 we do
41  compare against [15] that uses additional labeled data from ImageNet. Since we are the first to propose semisupervised
42  learning for adversarial robustness, there were no previous methods for using unlabeled data that we could directly
43  compare to. Nevertheless, in Appendix C.1 (described in lines 229–234), we compare our proposed method (RST) to a
44  state-of-the-art method for standard semisupervised learning (VAT), which we adapt to the robustness setting. In this
45  comparison, both methods use the same unlabeled data, and we find that RST offers significantly stronger performance.

46  *"Minors."* The reviewer raises questions about the relation between results in Table 1 and Figure 1b, both of which report
47  robustness to $\ell_\infty$ attacks but differ in a crucial aspect: Table 1 shows results for *heuristic defenses* tested against strong
48  gradient-based attacks (with $\epsilon = 8/255$), while Figure 1b compares methods with *certified (proven) robustness* against
49  *all attacks* with (with $\epsilon = 2/255$). State-of-the-art heuristic defenses utilize different algorithms than state-of-the-art
50  certified defenses and therefore we cite different works in each case. Moreover, since proving robustness to all attacks is
51  a harder problem than defending against particular attacks, the former comes at a cost to standard accuracy, explaining
52  the difference the reviewer points out.