

1 We thank all the reviewers for their valuable comments. We would be happy to revise our paper using their suggestions.

2 **Reviewer 1: Lipschitz constant and adversarial robustness.** While we agree that different approaches can be used
3 to improve adversarial robustness, optimizing the Lipschitz constant of a model is shown to tightly bound the *worst-case*
4 error due to bounded input perturbation (See Fig 1(e) below from our experiments; our SI and *Fazlyab et al., NeurIPS*
5 *2019*). The Lipschitz constraint is particularly effective at targeting precisely the regions of peak sensitivity of the
6 model, and our approach thereby provides provable guarantees for the trained model against adversarial perturbations.

7 We respectfully disagree with the statement that a “model can have large Lipschitz constant but can still be robust”. In
8 the provided example, if g achieves a higher performance than f , then the model g lacks robustness. In fact, a very
9 small perturbation (of the order of $1/1000$) can degrade the performance of g at all testing points close to the decision
10 boundaries (which are, in the absence of perturbation, correctly classified due to the high-frequency component). If g
11 achieves the same performance as f , then g contains an unnecessary high-frequency term. In fact, because our loss
12 function is strictly convex, it admits a unique minimizer, and high-frequency components are present in our models only
13 if they improve the overall performance (our Laplacian approach favors smooth solutions). This analysis is compatible
14 with the result that the Lipschitz constant provides a worst-case bound on the robustness of a model to perturbations.

15 Finally, despite its title, the reference cited by Reviewer 1 does not claim that bounding the Lipschitz constant is not a
16 good approach to obtain a robust model. Rather, it argues in favor of our work: it claims that the existing methods to
17 compute the Lipschitz constant had limitations, and that overcoming them would lead to adversarially robust models.

18 **Reviewers 1 and 4: Empirical validations, figures, and reproducibility.** We thank the reviewers for pointing out
19 these weaknesses. We recreated our figures (below; which now align with Reviewer 1’s intuition). We can now provide
20 more MNIST implementation details, upload our code to GitHub, and further demonstrate our robust training scheme.

21 **Reviewer 2: Confidence vs accuracy.** Confidence and accuracy are directly related (Figure 1(d)). Confidence is a
22 smoother metric for optimization and encoding the Lipschitz constraint, and is readily mapped onto one-hot vectors.

23 **Reviewer 2: Scalability.** We are able to successfully demonstrate our approach using the MNIST dataset. While
24 additional studies are needed to investigate scalability, our numerical studies show that our performance increases
25 rapidly with the number of vertices (model complexity), with an accuracy of 96% with only 10000 vertices (Figure 1(c)).
26 For comparison, our study shows that neural networks with the same number of parameters achieve only 68% accuracy.

27 **Reviewer 2: Dataset.** We choose the checkerboard dataset precisely because it requires a Lipschitz constant that tends
28 to infinity, so as to compare the performance over a broad range. We use the MNIST dataset for a more realistic study.

29 **Reviewer 3: Conclusion and novelty.** To the best of our knowledge, our results are the first to prove that a fundamental
30 tradeoff exists between Lipschitz constant and accuracy. Prior works (cited) provide empirical or less general results.
31 Our paper contributes both new results to the literature to address the immediate questions on tradeoffs, and a novel way
32 of looking at the learning problem (connections to PDE) that has broader scope and will pave the way for future works.

33 **Reviewer 3: Connections between problems and motivation.** The first problem generates models with minimal loss
34 and desired robustness. The second problem fixes the performance (loss) and improves the robustness of a model. While
35 the first problem yields a training scheme, the second is used to prove a fundamental tradeoff between performance and
36 robustness. We have drawn motivation for our work from prior works on adversarial attacks that have demonstrated a
37 glaring problem of robustness of machine learning models and a need for adversarially robust training.

38 **Reviewer 4: Discretized models and neural networks.** The discretized versions of the problems offer a way to
39 construct an implementable training scheme with robustness guarantees. Implementable models are obtained by various
40 discretizations of the original infinite-dimensional loss minimization problem. We have considered a graph discretization
41 that directly inherits the convexity of the infinite-dimensional problem. As the reviewer points out, other function
42 space parametrizations, such as neural networks, result in non-convex finite-dimensional optimization problems. Our
43 infinite-dimensional analysis holds regardless of parametrization and characterizes the solution that particular models
44 approximate, while our graph-based design offers a novel and alternative scheme to construct provably-robust models.

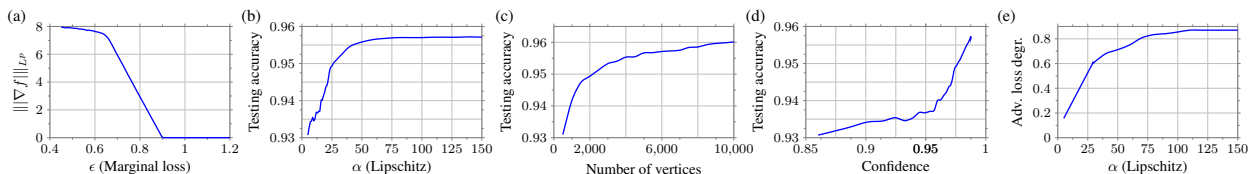


Figure 1: Figure (a) updates Fig. 2c in the paper. Figures (b)-(e) complement our study of the MNIST dataset by looking at the relations between accuracy and Lipschitz, accuracy and complexity, accuracy and confidence, adversarial robustness and Lipschitz.