1 We will fix all minor comments and typos without explicitly addressing them in the rebuttal.

2 **Response to Reviewer 1**:

3 *Practical Impact*: Our primary aim in this work is indeed theoretical. There has been substantial interest in the
4 theoretical understanding of adversarial robustness recently. Our work highlights the deficiencies in some of these
5 theoretical formulations (see also response to Reviewer 2 below), which we hope will lead to better theoretical models,
6 which in turn may lead to practical advances. Regarding an algorithm for monotone conjunctions in Theorem 10's
7 setting, the standard PAC learning algorithm for conjunctions suffices. An outline of this already appears in the
8 Appendix, but we will add a reference to it in the main paper.

9 *PAC Terminology*: We have assumed that readers will be familiar with standard terminology from PAC learning. Given
10 that many NeurIPS attendees may be unfamiliar with this terminology, we will add an appendix giving definitions that
11 we require and point readers to standard texts for further details.

12 *Non-trivial Class*: The definition of non-trivial class appears just before the statement of Theorem 5 (in lines 182-183).

13 *Undefined Algorithm*: The algorithm for *exact learning* monotone conjunctions using membership queries would be
14 considered folklore in the computational learning theory world; the key idea is that starting from the instance where all
15 bits are 1 (which is always a positive example), we can test whether each variable is in the target conjunction by setting
16 the corresponding bit to 0 and requesting the label. We will add this in the aforementioned new appendix.

17 *Finite Concept Classes*: Since (Thm. 11/Prop. 12) are primarily concerned with showing hardness of robust learning,
18 we don't think finite concept classes is a restriction. Please also see lines 70-90 for discussion regarding concept classes
19 defined over $\mathbb{R}^n$.

20 *Experiments*: We do not believe that *artificial* experiments will add to the value of the paper; that's not the main point
21 of the submission.

22 *Comparison to Prior Work/Contributions*: We will expand on the section in the paper, but we also refer to the review by
23 R3, which we believe very clearly summarizes our contributions.

24 **Response to Reviewer 2**: *Right Model*: We obviously disagree with the reviewer about this being a bad paper, but to a
25 great extent do agree with the reviewer about these being *unsuitable models* or *inadequate definitions* for adversarial
26 robustness. The point is that we *weren't* the inventors of these definitions (cf. [4, 5, 7, 8] for theory papers and others
27 more applied papers [A, B, C]). Our aim was precisely to show that once these definitions are *accepted*, even the most
28 elementary classes prove to be hard to robustly learn—and that proving computational hardness is much easier and
29 straightforward compared to the proofs that appeared in prior work.

30 Having criticized the definitions, we should acknowledge the contributions of prior work. Indeed, our initial aim was to
31 show positive results for at least some non-trivial classes under these definitions. It is clear that these definitions are in
32 many ways *natural* and *reasonable*, but when one looks at them under the lens of computational learning theory their
33 inadequacies surface immediately. We hope our work will highlight these issues and lead to future work (including
34 hopefully by us) that comes up with definitions that still (somewhat) retain the *simplicity* and *naturalness* of the current
35 definitions, while allowing one to separate non-trivial classes that are easy to robustly learn from those that are not!

36 *Computability/Halting Problem*: There is no connection to the halting problem, which is only one of the reasons why
37 uncomputability arises; there can be several others. The difficulty in this case is *enumeration* over (uncountably) infinite
38 sets. How would one compute the function, $\mathbf{1}(\exists y \in B_\rho(x).h(y) \neq c(y))$ in *finite time* even if one had black-box access
39 to evaluate $h$ and $c$ (the latter is not possible without membership queries)? Even under the assumption that the Turing
40 machine has the power to perform arithmetic over reals in unit time, the existential quantifier makes evaluating the
41 robust loss impossible! Even if the instance space is $\mathbb{Q}^n$, the decision problem for detecting an adversarial example
42 would be *recursively enumerable*, but not *recursive*. This problem disappears for finite instance spaces, but even there it
43 is not obvious how to evaluate this loss *without* membership queries. This is why one gets the separation for monotone
44 conjunctions depending on whether or not the learning algorithm has access to membership queries. In the case of
45 infinite instance spaces, we can't see a way to avoid the enumeration question without a strong inductive bias on $h$ and
46 $c$; in that case, properties of these functions, e.g. Lipschitzness, could be used to compute the loss in finite time.

47 **Response to Reviewer 3**:

48 We thank the reviewer for the comments and will obviously fix the typos they observed.

49 A. Cisse et al. Parseval Networks. ICML 2017.
50 B. Madry et al. Towards deep learning models that are resistant to adversarial attacks. ICLR 2018
51 C. Tramer et al. Ensemble Adversarial Training. ICLR 2018