1 We thank all reviewers for their feedback. We respond to the main concerns. We will also address all the minor points
2 raised in the final version of the paper.

3 **Motivation / Threat Model / Use Case (Reviewers 1, 2, 3)**   The main contribution of our paper is to provide a
4 framework that cleanly models privacy against bounded adversaries, and allows **quantitative calculation of privacy**
5 **loss** against particular adversarial classes. This was previously unknown (except for the special case of computationally
6 bounded adversaries). Quantification of privacy loss is important because it allows a systems-designer to precisely
7 determine the kind of privacy-accuracy tradeoff that is offered by a release.

8 One setting where quantifying privacy for different adversaries makes sense is when data sharing is coupled with data
9 usage contracts (as mentioned in Section 1 of the paper). For example, an instance of the Laplace mechanism might
10 offer $\epsilon = 1$ in general but $\epsilon = 1.25$ to a certain class of adversaries. Quantifying this tradeoff allows (a) better decisions
11 in cases where we expect the adversaries to be bounded in what they can do – for example, automated adversaries or
12 adversaries under a data-usage contract – and (b) better design of data-usage contracts – eg, if the loss against quadratic
13 functions is much higher than linear, then we can only allow for linear functions.

14 R3 astutely observes that problems may arise if the adversary does not obey the data usage contract, or its output is
15 viewed by someone else. In this case, we will sacrifice the improved privacy guarantee, but if we use a differentially
16 private mechanism, then we can fall back to the original (weaker) differential privacy guarantee that holds for all
17 adversaries.

18 We will add this discussion to the introduction in the final version of the paper.

19 **Reviewer 1**   "The only thing missing from the paper is a better sense of how adversary types can be linked to function
20 classes, perhaps through an extensive example."

21 A concrete example is an excellent suggestion; we will do this in the final version.

22 "Why couldn't the adversary just use simulation and rely on $\mathcal{H}$ instead? Why wouldn't the adversary have access to the
23 output of the cb-DP algorithm?"

24 To clarify, the adversary here can only compute certain functions (those in $\mathcal{H}$) on the output of the capacity bounded
25 differential privacy algorithm. The above discussion about data usage agreements is one setting where adversaries of
26 this form arise.

27 We will correct the error pointed out in the Appendix as well as add theorem numbers to our references where applicable.

28 **Reviewer 2**   "At this point it feels that the risks (in terms of possible privacy breach) of using this relaxation outweighs
29 the potential benefits."

30 We emphasize that the main contribution of the paper is quantifying the risk for capacity bounded adversaries. As
31 discussed above, under a mechanism (like a data usage agreement), one can get adversaries to only postprocess the
32 outputs of a mechanism using restricted function classes to get a tighter privacy analysis. If the adversary deviates, one
33 could still fallback to the general (and weaker) DP guarantee.

34 **Reviewer 3**   "The main problem here is that the definition - as the submission correctly observes - is not closed under
35 post-processing. In other words, once the (ML or contractually bound) adversary does its computation, _its_ output can
36 be observed and processed by someone else, without restrictions imposed by the definition."

37 Data use agreements usually restrict the outright release of data or its derivatives. This would ensure that the output of a
38 capacity bounded DP algorithm is not released to an adversary of a different class. That said, if this does happen in
39 violation of the data use agreement, we can still fall back on the (weaker) DP guarantee that applies to all adversaries.

40 We will qualify our statement on the Groce et. al. paper in Section 1, as was suggested. We will also add the suggested
41 details to the definitions section.