1 Thank you for your insightful feedback, corrections, and additional references that we will incorporate into our paper.
2 Below, we address several key points.

3 **All Reviewers (particularly R3):** "scope is a bit narrow" & "significance seems limited" - Based on consistent feedback
4 from the reviewers, we see that our presentation currently limits the perceived relevance and general importance of
5 our work. In our introduction, we will emphasize that our work is directly applicable to any scenario that requires
6 computing confidence intervals around importance sampling (IS) estimates. More broadly, we will also discuss that the
7 community is interested in our definition of safety [38] and its limitations [39], and IS [4 (in paper), 40, 43, 44]. Lastly,
8 we will mention that the $\alpha$-security formalization also pertains to high-confidence methods that do not use IS [31 (in
9 paper), 41, 42].

10 **R1:** "generalize to continuous MDPs" - Our work generalizes to continuous MDPs, but care must be taken to select $\pi_e$
11 such that IS weights are bounded. For example, the diabetes treatment simulation discussed in the paper has continuous
12 states and actions. We will make sure to discuss this extension in the paper.

13 "problems in applying Hoeffding's inequality" - In order to use Hoeffding's inequality, we assume that the IS weights
14 are bounded. Although WIS is biased, it works very well in practice. Consequently, we introduced the notion of
15 quasi-$\alpha$-security in Definition 2 to specifically allow for the analysis of WIS.

16 "if the attacker knows we are using Panacea" - The optimal attack does not change (lines 261–262), and therefore,
17 Panacea limits the damage incurred by the attacker. We will move this discussion outside of the proof block.

18 **R2**: "include studies on how easy it is to trick those algorithms too" - We are definitely interested in pursuing follow-up
19 directions to ensure security for model-based approaches, which we predict would be quite different and a significant
20 contribution on its own. In a future work section, we will include a discussion of similarities and additional challenges
21 that arise in that setting.

22 "the way $\pi_e$ is chosen" - We completely agree on the importance of how $\pi_e$ is chosen, even though the violation of
23 safety comes primarily from the safety test. Our current definition of security assumes that $\pi_e$ has lower performance
24 than $\pi_b$, but does not specify how often this occurs. Attacking the data used to select $\pi_e$ can increase this frequency.
25 Notice that attacking the data used to select $\pi_e$ alone would not cause the safety property to be violated. We will add a
26 discussion on this topic.

27 "whether the trajectory must still have been performed in the real environment" & "single out the few trajectories" -
28 Because the transition and reward functions are not known, one can not distinguish real and fake trajectories. Rare
29 events are critical to account for, and may look like fake trajectories. Perhaps impossible trajectories can be identified
30 using domain-specific knowledge, but that must be analyzed on a per-domain basis. We will mention this in the paper.

31 "use of any $\alpha \geq 1$ would be pointless" - We see that we did not provide sufficient discussion of Table 1, saying that the
32 behavior you note is what we aim to show! The middle column is usually $\geq 1$, indicating that standard methods can be
33 completely broken (make pessimal policies appear optimal) easily, as you described. However, the right column shows
34 values of $c$ that make Panacea $\alpha$-secure for *any* $\alpha \in [0, 1]$. E.g., plugging in $\alpha = 0.05$ gives Panacea a meaningful
35 security guarantee. Note that if $c \leq 0$, Panacea is not useful, but that the values of $c$ are positive and grow quickly as $n$
36 grows relative to $k$.

37 **R3:** "a worst-case stand-in" - When the stakes are high – for example, in the application of RL to sepsis treatment in
38 the intensive care unit, wherein training data is generated from hand-written doctors' notes – we do not want to assume
39 that the data contains only minor errors (such as patient height), but also major ones (such as wrong drug name).

40 **R4:** "an upper bound on the number of corrupted samples" - We will add a discussion of the many issues faced by
41 practitioners, including estimating the number of corrupt trajectories (perhaps based on known error rates in the data
42 processing pipeline of NLP and computer vision models) and selecting $\pi_e$. Panacea is only one piece of the puzzle, but
43 provides guarantees that are informative to practitioners.

44 **References:** [38] Ghavamzadeh, Mohammad et al. Safe policy improvement by minimizing robust baseline regret.
45 NeurIPS 2016; [39] Guo, Zhaohan et al. Using options and covariance testing for long horizon off-policy policy
46 evaluation. NeurIPS 2017; [40] Jiang, Nan et al. Doubly robust off-policy value evaluation for reinforcement learning.
47 ICML 2016; [41] Kuzborskij, Ilja et al. Confident Off-Policy Evaluation and Selection through Self-Normalized
48 Importance Weighting. arXiv preprint arXiv:2006.10460 2020; [42] Laroche, Romain et al. Safe policy improvement
49 with baseline bootstrapping. ICML 2019; [43] Liu, Qiang et al. Breaking the curse of horizon: Infinite-horizon
50 off-policy estimation. NeurIPS 2018; [44] Mandel, Travis et al. Offline policy evaluation across representations with
51 applications to educational games. AAMAS 2014.