1 **[General Response]** We thank all reviewers for the detailed and valuable feedback. We will fix the typos and improve
2 the draft carefully based on your comments.

3 *Q1.Limitation of the deterministic setting and extension to stochastic settings:* Our main contribution is designing
4 a simple value iteration style algorithm to get the interval estimation. For the sake of clarity, we choose to focus on
5 the deterministic transition and reward settings. One way to extend our results to stochastic settings is to follow the
6 Appendix C proof's technique in [1], by decomposing the empirical operator $\hat{\mathcal{B}}^\pi$ as $\mathcal{B}^\pi + (\hat{\mathcal{B}}^\pi - \mathcal{B}^\pi)$, where we can
7 bound the latter operator $(\hat{\mathcal{B}}^\pi - \mathcal{B}^\pi)$ via Rademacher complexity of the Lipschitz function class. To avoid digressing
8 from the primary focus of the current paper, we decide to leave them as our future work.

9 **[Reviewer #1]** We thank reviewer #1 for your valuable suggestions and detailed writing corrections.

10 *Q2.Misleading claim on non i.i.d. assumption:* Even when assuming deterministic transitions, typical approaches
11 based on concentration inequalities still require i.i.d. conditions on the transition $(s_i, a_i)$ pairs, or can only work on the
12 trajectory level as in [2] (which has a smaller effective sample size). We plan to have a new section to compare the
13 concentration inequality approach with our method side by side to further clarify the raised concerns.

14 **[Reviewer #2]** We thank reviewer #2 for your insightful questions and valuable experimental references.

15 *Q3.Empirical comparison to existing approaches:* We have a comparison with [2] in Table 1 of Appendix D. In general,
16 [2] views each trajectory as one sample while we view each transition pair as one sample. The result in Appendix
17 D shows that our method is better than that of [2] when the number of trajectories is small. Moreover, like other
18 trajectory-based IS methods, [2] also suffers from the curse of horizon. We will add more discussion in the revision.

19 *Q4. Reason to choose Lipschitz function class:* We had a brief discussion on this in line 128-129. We choose the
20 Lipschitz function class to make a good balance between expressiveness and tightness of the bounds. More specifically,
21 it includes a very rich set of functions that could cover the true value function with high probability, while allowing us
22 to get practical bounds with a simple algorithm. We will add more discussion in the final draft.

23 *Q5.Distance function and high dimensional state space:* We find out $L_2$ distance is enough for the low-dimension
24 environment. In high-dimension data, we may need to find a better distance measure to capture the underlying
25 low-dimension manifold of the data, which seems to be an exciting direction to explore. We will leave it as future work.

26 **[Reviewer #3]** We thank reviewer #3 for your valuable comments and suggestions.

27 *Q6.Quadratic dependency on sample size and random sample method seems not ideal:* We avoid the quadratic
28 dependency by adopting the random sub-sampling technique, which may sacrifice the tightness of the interval bounds
29 for reducing the computation burden. Moreover, as the sub-sampling bounds are still provably correct bounds of the
30 true $R^\pi$ (despite being less tight), the sub-sampling interval bounds can still guide the end-user for decision-making. In
31 real world applications, we can trade off the tightness and the computation complexity conditioning on the available
32 computation resource.

33 *Q7.Require knowledge of Lipschitz constant:* We agree that Lipschitz constant is crucial for the success of a valid
34 interval estimation. We emphasize and discuss it in section 4.2.

35 *Q8.Related literature on estimating distribution min/max form which can improve equation (14):* Thanks for pointing
36 out the reference. It sounds very interesting! Could you kindly send the references in the revised rebuttal?

37 **[Reviewer #4]** We thank reviewer #4 for your valuable comments and suggestions.

38 *Q9.The sample complexity appears to be exponential in the effective dimension:* We agree that the sample complexity
39 is exponential and the main reason to choose Lipschitz function class is because it strikes a good balance between
40 richness and simplicity. In line 180 we pointed out that: "it is possible to choose smaller space sets (such as RKHS) to
41 obtain smaller gaps, it would scarify other properties such as capacity and simplicity."

42 *Q10.Theorem 3.2 for the specified initial points:* We can start with an arbitrary initial point and still achieve linear
43 convergence as in Proposition 3.3, which means that when the algorithm converges we will get a pair of provably
44 correct bounds. However, with the specified initial point in Theorem 3.2, we can have a stronger guarantee on **anytime**
45 bounds, which means that whenever we stop the algorithm (before it converges), the upper and lower bounds we get is
46 guaranteed to include $R^\pi$. Moreover, we believe that the calculation of the initial point is not difficult (see Eq (12)).

47 **[References]**

48 [1] Mousavi, Li, Liu, Zhou. Black-box Off-policy Estimation for Infinite-Horizon Reinforcement Learning, ICLR 2020.

49 [2] Thomas, Theocharous, Ghavamzadeh. High-confidence off-policy evaluation. AAAI 2015