1 Thank you for your in-depth and constructive reviews, they will give us an excellent chance to further improve the paper.
2 We address the reviewers concerns individually below. We will not address all the typos and editing catches, but will fix
3 all of these in the final draft.

**Reviewer 1:**

5 - The goal here is to compare model-free methods, augmented with a buffer. A careful comparison between model-free
6 and model-based approaches for OPE would be interesting and extremely valuable for the community, but is beyond the
7 scope of this paper.

8 - It is a good suggestion to provide intuition for the proofs. We will include such a discussion in the camera-ready,
9 which allows for an additional page.

10 Minor Concerns: We will address these concerns in the revision. We will improve notation consistency and clarity and
11 include diagrams of the maze in the appendix.

**Reviewer 2:**

13 - SIR is a general strategy; in fact, there are a number of similar approaches with different names [Smith and Gelfand,
14 1992]. The main novelty here is investigating its use in RL, where the online setting requires us to consider a moving
15 window dataset—rather than a fixed batch.

16 - The theoretical comparison of the bias of IR and WIS-Optimal is natural, because we show they are equal. IS is
17 unbiased, so that comparison is not interesting. In practice, though, we cannot actually use WIS-Optimal, as it is a full
18 batch approach. Empirically, then, it makes sense to compare to other mini-batch methods, like IS. We did not compare
19 to WIS-minibatch due to the poor empirical performance, likely due to the additional bias of that estimator.

20 - $\bar{\rho} \approx \mathbb{E}[\rho(a|s)] = \mathbb{E}[\frac{\pi(a|s)}{\mu(a|s)}] = \sum_{s,a} \frac{\pi(a|s)}{\mu(a|s)} \mu(a|s) d_\mu(s) = 1$.

21 - Assumption 1 is common for analyzing OPE estimators. The idea is that we are effectively sampling from the stationary
22 distribution, even though we know we are in Markov settings. An important next step is to consider alternative noise
23 assumptions in sampled data.

24 - The result in line 203 is that we can directly use the prior results to look at the expected difference in variances over
25 many buffers (i.e. these statements say our result holds across buffers of smaller sizes).

26 - Because $\bar{\rho} \approx 1$. When $\rho$ is lower than the average it will make the rhs a large number, but when $\rho$ is greater than the
27 average we expect it to lower the rhs of the equation. As learning progresses, we expect the samples w/ high $\rho$ to learn
28 more quickly (thus having lower error). 'mean' is the correct one.

29 - We also looked at "softer" target policies, where similar conclusions can be drawn (see appendix). All the results
30 presented in the appendix are qualitatively similar.

31 - The parameter sensitivity provides more information, because it gives some sense of how these might perform in
32 practice for realistically chosen parameters, rather than optimal parameters. We will include the learning curves in the
33 appendix for completeness.

34 - MARE is Mean Absolute Return Error. We use MARE when it is not tractable to compute the value function using
35 dynamic programming or analytically (and otherwise MAVE).

**Reviewer 3:**

37 For O1, you are correct in your understanding. We will use some of the additional space in the camera-ready to include
38 a brief discussion on O1 and O2.

39 Q1: The Steps corresponds to the Number of Interactions with the environment. The agent can update more or less
40 frequently than every step. The Number of Updates for Figure 1(a) is once every 16 Steps. We also show performance
41 for different Update frequencies in Figure 1(b).

42 Q3: Because the experiment is run 100 times, like in all the plots in figure 3, the error bars are not visible. The parameter
43 sensitivity plots could provide some information about variance of the updates. If the variance of the updates is higher,
44 we expect the magnitude of the largest updates to also be higher. This means a lower step size is needed to prevent
45 divergence. A wider trough of the sensitivity curve could reflect lower variance in the updates, though as acknowledged
46 in the paper, this is very much a proxy and we cannot make any strong conclusions based on it.

47 C1(f): See point 3 for Reviewer 2.

48 Q2, C1(a-e), C2, and C3: we will take all these points into consideration, and work to maximize clarity in the final
49 revision. L287 - You are correct, this should be figure 3.