1. We thank the reviewers for their valuable feedback. Please see the following for the response, and we will make
2. corresponding modifications in the revised paper.

3. R2 and R3 had questions about the **time complexity** of our method. The only additional cost of our method over
4. standard RL is running graph search, which must be done once per evaluation episode. On an 8-core desktop with
5. 32GB of RAM, it takes 65 seconds to estimate the $\mathcal{O}(|V|^2)$ pairwise distances between the 1000 observations in the
6. replay buffer. As noted in Appendix A, this computation can be amortized across many goal-reaching tasks. Given
7. that pairwise distances, it takes 0.029 seconds to compute the shortest path (using Dijkstra's Algorithm, which runs in
8. $\mathcal{O}(|V|^2)$). In terms of **sample complexity**, we found that SoRB required more samples to converge than SPTM, likely
9. because dynamic programming (the optimization problem for SoRB) is more challenging than supervised learning
10. (the optimization problem for SPTM). Lastly, we agree with R2 that the construction of "good" replay buffers is an
11. important problem to explore.

### Reviewer 1:

13. • $d_{s \to g} > $ MAXDIST – For most distant goals, the shortest path does not go through the waypoint, so the distance
14. directly to the goal is almost always shorter. However, always conditioning on the goal would be problematic, because
15. the goal-conditioned policy often fails to reach distant goals directly (See Fig. 6). We will clarify this in Section 2.3.
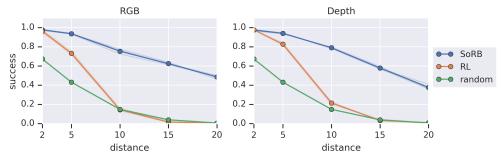
### Reviewer 2:

17. • **Replay buffer**: For the generalization experiments, we initially seeded the replay buffer with 1,000 observations
18. taken from random states and actions in the new environment. We will clarify this in Alg. 1.

19. • **Where should SoRB excel?** We expect SoRB to outperform existing RL approaches on long-horizon tasks,
20. especially those with high-dimensional inputs. The Atari games fit that mould, as does StarCraft II. Many of the
21. Mujoco tasks in the OpenAI Gym don't require long-horizon reasoning, so we expect little benefit from using SoRB
22. for environments like Hopper and HalfCheetah. We will add this discussion to Section 6.

23. • **Comparison with state representations** – We agree that learning compact state representations might boost the
24. performance of the underlying goal-conditioned agent. While we found Value Iteration Networks (which learn
25. a representation based on the values of states) performed poorly (see Fig. 6), it is plausible that other methods
26. (e.g., [14]) would work better. We will include a comparison to [14] for the final version.

### Reviewer 3:

28. • **RL and Planning, PRM-RL** – We discussed some prior work that combines RL and planning, including PRM-RL,
29. on L165 - 167. The key difference between PRM-RL and SoRB is graph construction. Whereas SoRB simply uses
30. the value function to predict the distance between two states, PRM-RL requires evaluating the policy in the real-world
31. to determine this distance.[12, Section III.B]

32. • **Why does Q-learning fail to generalize?** – While, in theory, pure Q-learning should be able to solve any task, in
33. practice the optimization problem is quite challenging, and only succeeds in learning to reach nearby goals. Our
34. method generalizes better than Q-learning because, while both are solving a dynamic programming problem at the
35. high-level, planning provides an exactly solution, while Q-learning gives an approximate solution.

36. • We will include a comparison to a hierarchical RL method (e.g., [36]) in the final version.

37. We ran an additional **generalization experiment**. The only difference from Fig. 9 in Section 5.5 is that we now
38. normalize input images to be in [0, 1] (i.e., divide RGB by 255). We found that normalization substantially improved
39. the generalization results for RGB images, reaching almost 80% of goals 10 steps away, while goal-conditioned RL
40. reaches less than 20% of these goals. Transparent lines correspond to average success rate across 22 held-out houses for
each of three random seeds.



41.