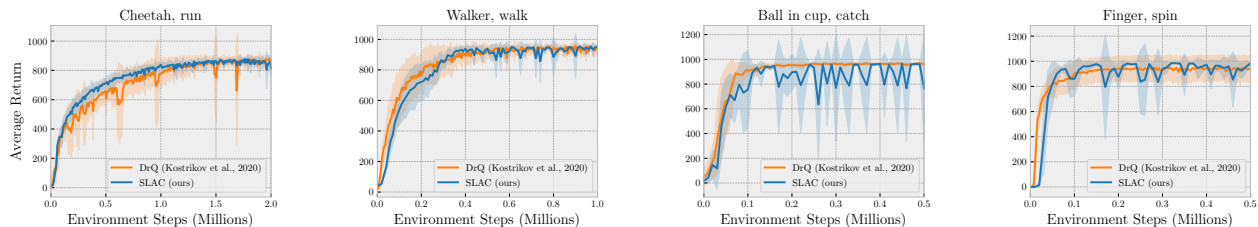


1 We thank the reviewers and are happy they found the paper is well-written (R1, R2, R3, R4), the method is theoretically  
2 sound (R1, R2), the experiments are appropriate and comprehensive (R2, R3, R4), the results are convincing (R1, R3,  
3 R4), and the ablation studies are “tremendously useful” and helpful for making design choices (R1, R2, R3, R4). We  
4 address individual questions below.

5 **R1: POMDP claims.** We’ll update the paper to stress that our method is not equipped to solve the POMDP, e.g.  
6 reducing uncertainty. It was not our intent to claim that it was. Although we use the POMDP formalism to derive our  
7 algorithm, we agree with R1 that our experiments do not address the capability of our method to solve POMDPs.

8 **R1: SOTA claims.** We’ll remove the SOTA claims in light of the recent works CURL, RAD, and DrQ [1]. We strongly  
9 agree that this is distracting and less important than the content of the paper, and we had not intended to present  
10 that as the main contribution. Regarding the results: to our knowledge, DrQ (which came out on 28 April 2020, and  
11 would be reasonably considered concurrent with NeurIPS submissions) is the best among these methods, and SLAC  
12 is comparable to DrQ in the 4 DM control tasks, as shown in Figure 1. Note that the DrQ and CURL papers report  
13 outdated and lower-performing results for our method, compared to the results from this submission. That said, we will  
14 remove claims about state of the art results, and add results for all methods.



**Figure 1:** SLAC (ours) achieves comparable performance as DrQ (Kostrikov et al., 2020 [1]) in the 4 DM control tasks. Note that the plots include the initial exploration steps (10k for SLAC and 1k for DrQ), in which data is collected from a uniformly random policy.

15 **R2: Divergence and instability.** We hypothesize that it’s caused in part by the constantly changing latent space of the  
16 model and the model overfitting to the most recent data (the replay buffer contains the 100 most recent episodes). We  
17 believe that our method could benefit from early stopping or a schedule for model updates (e.g. having a "target" model,  
18 similar to target networks in RL). This merits further investigation, which we leave for future work.

19 **R2: Task-dependent representations.** The latent representation does receive some task-dependent information (in  
20 the form of reward prediction), although we acknowledge it could further incorporate more feedback from the value  
21 functions and policy. We chose to keep them separate but we recognize that it’s worth investigating further.

22 **R3: "The two main components of the algorithm can be drawn from existing methods in a straightforward  
23 way."** We acknowledge that the individual components are based on previous work, and we will discuss this more  
24 clearly in the paper. Although the resulting algorithm is simple, we believe that formally establishing a link between  
25 them is important to help us understand the implications of certain design choices. E.g., one interesting finding is that,  
26 as indicated by Equation (21) in the Appendix, the optimal policy of the SLAC objective is optimal with respect to  
27 the expectation over the belief of the Q value of the learned MDP. This is equivalent to the Q-MDP heuristic, which  
28 "amounts to assuming that any uncertainty in the agent’s current belief state will be gone after the next action" (Littman  
29 et al., 1995 [2]). We’ll update the paper to emphasize these observations.

30 **R3: "Overly optimistic / risk seeking."** We believe we don’t have that same issue since the agent isn’t allowed to  
31 control the dynamics of the *future* (i.e., the policy objective in Equation (10) updates only the policy parameters).  
32 Furthermore, the history-dependent policies mitigate a similar issue, except that it’s caused by something different:  
33 When executing a latent-conditioned policy, the latent distribution is available but the specific latent state is unknown.  
34 Any attempt to choose any latent (e.g. the mode or a random sample) results in an overly confident agent since the  
35 policy assumes that the given latent is correct. By conditioning the policy in the observations, we instead obtain the  
36 Q-MDP heuristic (see paragraph above). We acknowledge that our empirical evaluation doesn’t exemplify this. This  
37 issue doesn’t arise if the belief is narrow. We’ll update the paper to clarify this.

38 **R3: Correctness and clarity.** We’ll update the paper with clarifications on the points brought up by R3. Most  
39 importantly, the true state of the environment is never available.

40 **R1: Reference in the related work:** We’ll add the missing reference to the paper.

41 [1] Kostrikov et al. Image Augmentation Is All You Need: Regularizing Deep Reinforcement Learning from Pixels. arXiv:2004.13649, 2020.

42 [2] Littman et al. Learning policies for partially observable environments: Scaling up. Machine Learning Proceedings, 1995.