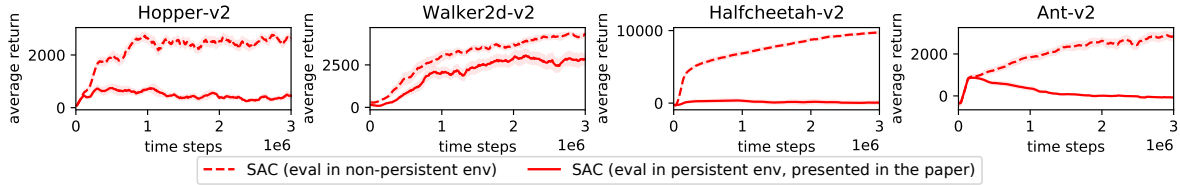
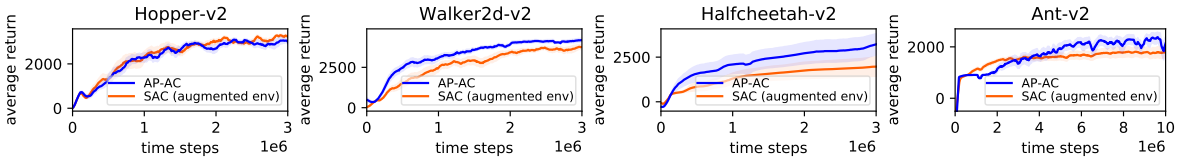


1 We thank all the reviewers for their valuable comments. We will reflect the comments in the final version of the paper.

2 **[R2: results on SAC]** Please let us correct the reviewer’s misunderstanding of our experimental results on SAC. In
 3 Figure 5 of the paper, the baseline SAC agent is trained on the standard **non-persistent** environments while being
 4 evaluated on **c-persistent** environments where the action-persistence is enforced. As shown in the following figure¹,
 5 the performance of SAC consistently improves in the non-persistent environment that the agent is trained on, but its
 6 naïve projection into a *c*-persistent policy completely fails since the agent *never* considers the action-persistence during
 7 training. We would like to emphasize this is natural and **not** a broken result.



8 **[R3: alternative baseline]** Your understanding is correct. The environment modification method, which includes the
 9 ‘last action’ and $t \bmod L$ in the augmented state, will also enjoy a linear complexity with respect to L . However, this
 10 alternative baseline still has some drawbacks compared to our proposed method. First, it is unable to exploit every
 11 transition sample to update every timestep’s actor and critic, while AP-AC is capable of doing it in Eq. (9-10). Second,
 12 there exists a redundancy in the representation of Q -function, i.e. $Q(t, \bar{a}_{last}, s, a)$ is not succinct compared to ours of
 13 $Q(t, s, a)$. This incurs a factor of $|A|^2$ increase in time complexity of policy evaluation in tabular FA-MDPs. Still, we
 14 conducted additional experiments that compare the proposed baseline (i.e. training vanilla SAC agent in the augmented
 15 environment that includes \bar{a}_{last} and $t \bmod L$ in the observation) to AP-AC. As the following figure demonstrates,
 AP-AC still performs better than or on par with the alternative baseline.

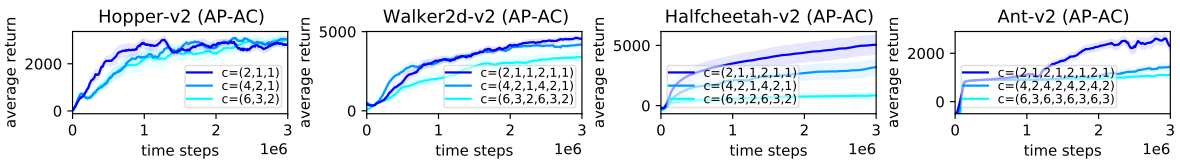


16 **[R3: experiments on standard benchmarks]** To the best of our knowledge, there is no standard benchmark to evaluate
 17 RL algorithms with multiple action persistence. Adopting our method to real-world tasks remains as future work.

18 **[R3,R4: related works]** Thanks for your suggestion. If every action variable’s control frequency is same, *c*-persistent
 19 action can be understood as a particular instance of an option in semi-MDP framework, which always lasts c time steps.
 20 We will add more discussions on related works such as semi-MDPs, temporally abstract actions, and action-repeats in
 21 the final version of the paper.

22 **[R4: comparison to Persistent FQI (PFQI)]** Our algorithm, AP-AC, is a non-trivial extension of PFQI. First, PFQI
 23 is only applicable to *single* action-persistence and finite action space, while AP-AC can deal with arbitrarily *multiple*
 24 action-persistence and both finite and continuous action spaces. Second, PFQI maintains only one Q -function and
 25 performs Bellman optimality backup followed by action-persistence backup k times. Each optimality (or persistence)
 26 backup operation requires to solve a regression problem until convergence. As a consequence, it can only work in
 27 the batch RL (a.k.a. offline RL) setting. Also, as long as PFQI maintains single Q -function, its extension to online
 28 RL is not straightforward. In contrast, AP-AC uses L actors and L critics and simultaneously updates all of them via
 29 exploiting their recursive relationship, which enables online learning of the agent as well.

30 **[R4: more comprehensive experiments]** The goal of this work is to provide an efficient solution method for the *given*
 31 action-persistence c , not finding a proper c to speed up learning. Still, we conducted additional experiments to present
 32 the effects on the resulting policy of varying c . As the following figure shows, larger action persistence yields more
 33 degradation of asymptotic performance due to a limited degree of freedom of control. AP-AC consistently works well
 34 for various c ’s. We will also include qualitative examples (e.g. videos) in the supplementary material of the final paper.



35 We will also place the Broader Impact section into the main text, which is currently in the supplementary material.

¹A performance gap exists compared to those reported in the original SAC paper, due to usage of different hyperparameters such as the number of hidden units per layer, i.e. 100 (ours) / 256 (original SAC paper).