

1 We thank the reviewers for their thoughtful feedback and suggested improvements. We agree that our proposed e-stop
 2 framework is a simple and attractive mechanism for incorporating state-only observations from an expert demonstrator,
 3 or manual interventions from a supervisor. Here, we provide clarifications and additional results which address the
 4 reviewer concerns necessary to further increase their scores.

5 **Additional empirical analysis** The central feedback from all reviewers was the need for experiments in a continuous
 6 MDP, in addition to our previously presented tabular results. We present new results on an inverted pendulum
 7 environment with continuous states and actions which will be included in the manuscript. We trained an agent using
 8 deep deterministic policy gradient (DDPG), with an actor and critic using two and three hidden layers of size 64,
 9 respectively. We collected 500 demonstration trajectories from a near-optimal expert policy and constructed a support
 10 superset from a rectangular hull of these samples. Fig. 1 compares DDPG with and without an e-stop mechanism, and
 11 shows roughly an order of magnitude improvement in sample complexity on a hold-out set of random seeds.

12 We provide further fine-grained analysis by controlling for several other parameters in the grid-world environment. To
 13 address [R1], we present results in Fig. 2 over a range of expert demonstrations used to construct the support superset $\hat{\mathcal{S}}$.
 14 Even after decreasing the number of demonstrations from 1000 to 5, the cumulative reward of the agent trained via
 15 our e-stop mechanism decreased by less than 3%. Note that a bound relating the number of expert trajectories to the
 16 suboptimality of our method was proved in Theorem 5.1. To address [R3], in Fig. 3, we controlled for the quality of the
 17 demonstrator policy by injecting random noise into the optimal policy’s Q -function. As predicted by Theorem 5.1, the
 18 learned e-stop policy has bounded regret with respect to the expert demonstrations, but in practice typically exceeded
 19 the performance of the expert.

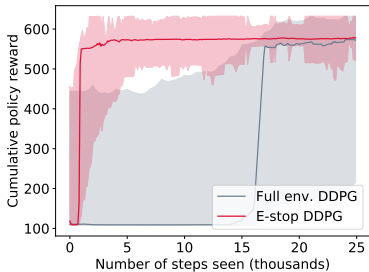


Figure 1: Pendulum environment.

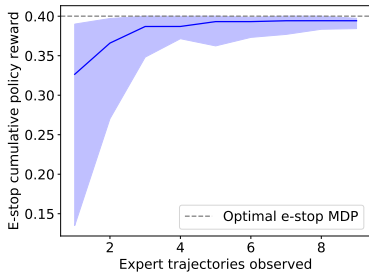


Figure 2: Number of observations.

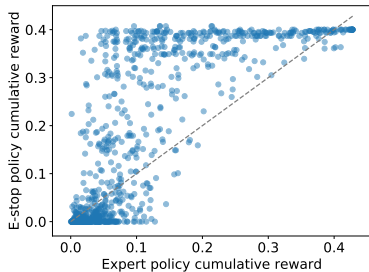


Figure 3: Expert quality.

20 **Relationship to existing work** We emphasize that our method applies to the reinforcement learning with expert
 21 observations (RLEO) setting, where only the states visited by the expert are observed. This precludes methods for
 22 the reinforcement learning with expert demonstrations (RLED) setting, where both the expert’s states and actions are
 23 provided. To answer [R2], it is not immediately clear how DQfD (Hester et al.) could be extended to the RLEO setting,
 24 as their method requires action observations to perform Bellman updates. However, it is possible to extend our method
 25 to RLED by constructing a support superset $\hat{\mathcal{S}}$ based on state-action pairs. Thus, our method and many RLED methods
 26 are complimentary. For example, DQfD would allow pre-training the policy from the state-action demonstrations,
 27 whereas ours reduces exploration during the on-policy learning phase. Similarly, our work can be related to Fujimoto et
 28 al. (as suggested by [R1]) by using a state-action superset and off-policy training, where states outside of $\hat{\mathcal{S}}$ are never
 29 selected in the Bellman update rule due to the maximum terminal e-stop penalty.

30 To answer [R3], the method in Eisenbach et al. is also complimentary to our work, and does not permit a direct
 31 comparison. To elaborate, they consider the RL setting where a second “soft reset” policy is learned in addition to the
 32 standard hard reset. The soft reset policy prevents the agent from entering nearly non-reversible states and returns the
 33 agent to an initial state. Hard resets are required whenever the soft reset policy fails to terminate in a manually defined
 34 set of safe states $\mathcal{S}_{\text{reset}}$. Our method can be seen as learning $\mathcal{S}_{\text{reset}}$ from observation (in the absence of a soft reset policy,
 35 we also use $\mathcal{S}_{\text{reset}}$ during the forward policy roll-outs). Their method trades off hard resets for soft resets, whereas ours
 36 learns when to perform the hard resets.

37 “Shielding” (Alshiekh et al.), as suggested by [R2], uses a manually specified safety constraint and a course, conservative
 38 abstraction of the dynamics to prevent an agent from violating the safety constraint. Our expert observations can be
 39 seen as inducing a safe region (the support superset), although our e-stop mechanism is model-free and uses a terminal
 40 reset penalty for actions which would have exited the safe region.

41 **Other minor improvements** We will make all typographical and clarity changes suggested by the reviewers, including
 42 using a more informative title.

43 We believe these changes, and in particular the additional results in continuous MDPs, address all reviewer concerns.