

772 **Generalization to Real Robots**

773 We were unable to validate our method on real robots due to the current COVID-19 pandemic. We  
774 believe our proposed method, LOOP, will also generalize to real-world robots for the following  
775 reasons.

776 **LOOP achieves better simulation results than prior work which has been demonstrated to**  
777 **work on real robots.**

778 **More sample efficient than SAC:** SAC [5] has been demonstrated to work well on real robots and is  
779 widely used in a number of robotics application [2, 60, 61]. We build upon SAC without making any  
780 algorithmic assumptions or having additional requirements (such as using a ground truth dynamics  
781 model). Also, Figure 3 shows that LOOP is consistently more sample efficient than SAC. This implies  
782 that it requires less deployment time of robot in real world and thus decreases wear-and-tear to the  
783 hardware. Hence we believe LOOP should generalize to real robots with the added improvement in  
784 sample efficiency.

785 **Less runtime than PDDM:** PDDM [21] was demonstrated to work on physical robots. In terms of  
786 trajectory optimization, LOOP optimizes for a shorter horizon utilizing its terminal value function and  
787 can achieve a better runtime performance on the real robot. In addition, PDDM also requires more  
788 hyperparameter tuning than LOOP as it sets different hyperparameters for different environments.

789 **LOOP is evaluated over a wide range of simulated tasks and RL domains.**

790 LOOP outperform the baselines across a range of simulated tasks consisting of Locomotion, Manipu-  
791 lation and Navigation experiments. Section 6 shows that LOOP shows improved sample efficiency  
792 in both the Online RL and Safe RL domain, making it a preferable choice to deploy with physical  
793 robots. Offline RL is an important problem in robotics as it allows for learning behaviors from static,  
794 previously collected datasets. LOOP shows versatility by improving over the performance of the  
795 parameterized actor in Offline RL as well.

796 **LOOP uses more interpretable policies that are more amenable to safety increasing its utility**  
797 **in real world applications**

798 LOOP uses semiparametric policies in the form of H-step lookahead. This way of online planning  
799 allows to interpret and reason about the outcome of the policy, while giving us flexibility to incorporate  
800 constraints during deployment. Incorporating constraints during deployment is important in physical  
801 robots where constraints could be safety considerations (possibly non-stationary). Online planning  
802 also allow LOOP to adapt faster to changing dynamics due to wear and tear of robot or any other  
803 phenomena.

804 **LOOP requires less hyperparameter-tuning for each task compared to model-based**  
805 **alternatives.**

806 Tuning hyperparameters in real world requires multiple environment runs for the experiments and  
807 is a non-trivial task. LOOP for Online RL works with a single set of hyperparameters across all  
808 environments and did not require much hyperparameter tuning. In contrast, other ways of utilizing  
809 model in off-policy methods like MBPO relies on hyperparameters specifically tuned for each  
810 environment as can be seen in Appendix C from [13]. The ease of hyperparameter tuning in LOOP  
811 will likely make it easier to be applied to real-robots compared to the previous works.