
Benchmarking Safe Policy Optimization for Constrained Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

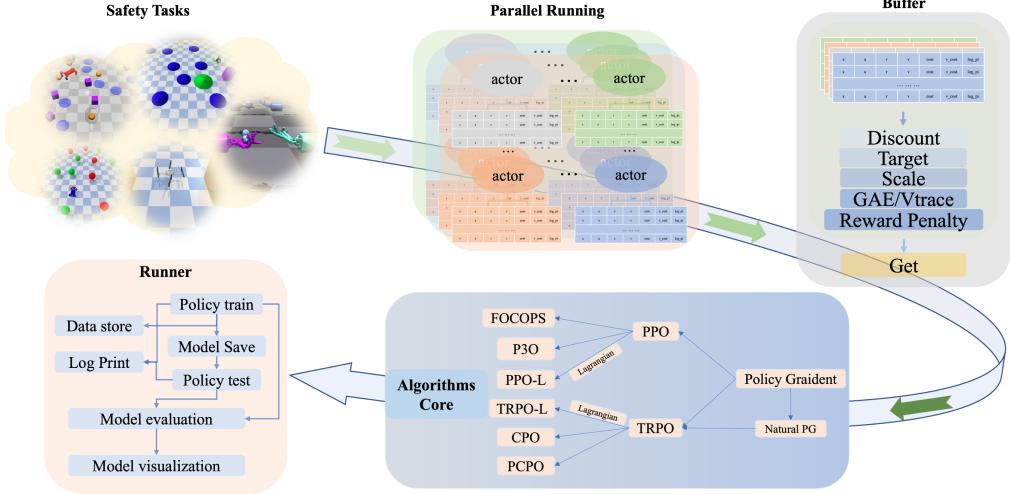
Safe reinforcement learning (safe RL) tackles decision making problems with safety constraints. Despite the influx of attention in this field, there is a lack of commonly recognized safe RL benchmark. This is partly because the code of many safe RL methods is unavailable and also because new methods often come with new testing tasks. As a result, researchers suffer from incorrect implementations, unfair comparisons, and misleading conclusions. In this study, we offer a solid safe RL benchmark—(SafePO)—which benchmarks popular safe policy learning algorithms across a list of common environments. Specifically, we start by standardizing the problem of safe exploration via solving constrained Markov decision processes (CMDP). Then, we provide implementations for CMDP solutions, covering both constrained policy optimization type methods and Lagrangian type methods. Our implementations in SafePO are highly efficient in the sense that learners can collect samples in parallel and synchronize their policy gradients on different physical CPU cores. We test them on four types of safety-aware robot learning tasks. Based on the benchmark results, we derive new insights by disclosing the interplay of different attributes on safety performance and illustrating the difficulty of safety learning on the sparse cost. Furthermore, to consider the safe RL problem in multi-agent settings, we introduce new tasks based on the DexterousHands environment and report comparison results for both single-agent and multi-agent safe RL algorithms. Our project is released at: <https://github.com/PKU-MARL/Safe-Policy-Optimization>.

1 Introduction

Safety policy learning is critical in real-world reinforcement learning applications, and dangerous decisions are undesirable. For example, a robot agent should avoid playing actions that irrevocably harm its hardware [43, 21]. Given its importance, recently the community actively consider safe policy learning (e.g.,[2, 47, 28, 59, 27, 54, 37]). However, most existing work only focuses on algorithm design, and it still lacks an open standardized algorithm benchmark for empirical study and a comprehensive comparison for the widely used safety tasks. This results in difficulty reproducing the experimental reports, unfair comparisons, and unclear understanding existing algorithms.

To address the undesirable research status of safe reinforcement learning, we open an efficient safe policy optimization benchmark and provide some empirical insights for safe reinforcement learning.

Standardizing Safety Learning Tasks and Algorithms. We propose the safe policy optimization (SafePO), which benchmarks typical and widely used safety tasks and learning algorithms. This benchmark collects significant diversity in the safety tasks, including Robots with Limited Velocity, Robots with Circle, Robots with Goal, and Safety DexterousHands. For the algorithmic part, this



(a) Overview of safe policy optimization

Figure 1: The proposed safe policy optimization (SafePO) provides a rich environment interface, which can be used for different algorithms to verify performance on different tasks. SafePO contains 6 baseline algorithms includes PPO-L, TRPO-L, P3O [60], PCPO [57], FOCOPS [61], and CPO [1]. Particularly, the SafePO collects samples in parallel and synchronizes the policy gradients on different physical cores of the CPU, which makes SafePO efficient to train a tough safety learning task. According to the algorithm configuration, the framework performs customized preprocessing on the collected data to update the actors and critics of different algorithms. During the algorithm training process, we also provide a large amount of visual information to evaluate the algorithms.

36 benchmarks includes CPO [1], PCPO [57], PPO-L, TRPO-L, FOCOPS [61] and P3O [60]. The
 37 proposed SafePO standardizes the safe policy optimization via the main formalism with respect to
 38 the constrained Markov decision process [43], which provides a training framework from compiling
 39 algorithms to the evaluation of algorithms. The SafePO collects samples in parallel and synchronizes
 40 the policy gradients on different physical cores of the CPU, which makes SafePO efficient for training
 41 a tough safety learning task. Besides, the proposed SafePO is a scalable software package that can be
 42 extended to new safety tasks at scale, which makes it an accessible way to promote the reproducibility
 43 of the latest safety tasks.

44 **Empirical Explanation for Safety Learning.** We consider the key factors that influence the perfor-
 45 mance of the safety learning, where we provide a comprehensive empirical comparison to benchmark
 46 algorithms with 33 different safety robot tasks setting. Our empirical observations include:

- 47 1) The first observation is that there is no unambiguous winner: no single algorithm performs best
 48 among all the tasks for reward performance and cost limit in all safety tasks. But in general, the
 49 approaches implemented by a deep neural network (e.g., P3O) achieve a better performance than the
 50 convex approximation method (e.g., CPO); such an empirical observation is sharper when the task is
 51 more complex (e.g., Safety DexterousHands, see Figure 6).
- 52 2) We observe that the safety task of robots with limited velocity is an *ill defined task*, where a
 53 contradictory phenomenon occurs between reward improvement and constrained limited cost for the
 54 robots like Ant-v3, Hopper-v3, HalfCheetah-v3, and Walker2d-v3, see Figure 3, and 7.
- 55 3) We observe that the Lagrange parameter plays an important role in learning a safe policy for the
 56 agent, see Figure 4. If the policy violates the cost limit, the Lagrange parameter will be active and
 57 help the agent update the policy alongs with a safe direction. Otherwise, the Lagrange parameter
 58 keeps calm, which allows the agent improves the reward performance within the safety region.
- 59 4) We observe that it is very challenging for the agent to learn a safe policy if the cost is too sparse.
 60 The empirical result shows that: the agent learns little to improve the reward performance and always
 61 prefers to play an unsafe policy. We suspect that sparse cost implies a small cost limit, which causes
 62 a limited safe region. Thus, it is very difficult for the agent to search for a safe policy.

63 **Safety DexterousHands Benchmark.** We propose the safety DexterousHands benchmark, a new
 64 and challenging robot task for measuring safe policy optimization algorithms. Dexterous manipulation
 65 with hands to complete the complex task is one of the most challenging problems in the robotics [38],
 66 the proposed Safety DexterousHands tasks provide the way to test the capacity of the algorithm to
 67 manipulate the tough tasks. For a deeper comparison to the safe policy optimization algorithms, we
 68 present the performance with two settings: single-agent and multi-agent.

69 **2 Related Work**

70 **Safety Algorithms.** The work [42, 4, 35] provides concrete problems in AI safety, for example,
 71 avoiding reward hacking and safe exploration. This benchmark mainly consider the safety learning
 72 via the formulation of constrained Markov decision process [3], which is an active community in safe
 73 reinforcement learning (e.g., safe exploration [40, 8, 31], risk criteria [11, 12, 44, 45], or primal-dual
 74 method [14, 49, 41, 33]). For a comprehensive survey, see recent works [29, 10].

75 **Reinforcement Learning Benchmark.** Arcade Learning Environment (ALE) [7] provides the
 76 widely used Atari games with score-based reward functions. Open AI Gym [9] includes wide
 77 variety of classical control tasks: MuJoCo [50], Atari games high-dimensional continuous control.
 78 The work [48, 20, 18] also proposes continuous robotic control tasks. Deepmind Lab [6] provides
 79 a first-person 3D game platform designed for the research and development of general artificial
 80 intelligence and machine learning systems. CoinRun [15] is a procedurally generated environment,
 81 and [15] shows the generalization performance during the training process. Surreal [23] is an open-
 82 source reinforcement learning framework and robot manipulation benchmark. Ciatech OPE [51]
 83 provides eight environments for the empirical study of off-policy policy evaluation. Pybullet [18]
 84 provides a framework that combines diverse learning paradigms (such as imitation and reinforcement
 85 learning). Assistive gym [22] is a physics simulation framework for assistive robotics. Brax[24] is a
 86 differentiable physics engine for large-scale rigid body simulation.

87 **Safety Benchmark.** AI Safety Gridworlds [35] provides a framework of safety reinforcement
 88 learning tasks with respect to grid world. The work [43] proposes Safety Gym that contains three
 89 robots for the empirical study of safe exploration in deep policy optimization. Bullet-safety-gym [27]
 90 provide a framework for safety specifications of Bullet [16] in constrained reinforcement learning
 91 problems. Safe-control-gym [59] provide an implementation for three dynamics. The work [29]
 92 considers safe multi-agent reinforcement learning benchmark.

93 **3 Safe Reinforcement Learning**

94 The paper formulates safe reinforcement learning (RL) as a constrained Markov decision process
 95 (CMDP) [3], which is defined as $(\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \rho_0, \gamma, \mathcal{C})$. Here \mathcal{S} is state space, \mathcal{A} is action space.
 96 $\mathbb{P}(s' | s, a)$ state transition probability from s to s' after playing a . $r(s'|s, a)$ denotes the reward that the
 97 agent observes when state transition from s to s' after it plays a . $\rho_0(\cdot) : \mathcal{S} \rightarrow [0, 1]$ is the initial state
 98 distribution and the discounter factor $\gamma \in (0, 1)$. The constraint set $\mathcal{C} = \{(c, b)\}$, c is cost functions:
 99 $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$, and b is cost limit. A stationary policy π is a probability distribution defined on
 100 $\mathcal{S} \times \mathcal{A}$, we use Π to denote the set of all stationary policies. Let $\tau = \{s_t, a_t, r_{t+1}, c_{t+1}\}_{t \geq 0} \sim \pi$ be a
 101 trajectory generated by π , where $s_0 \sim \rho_0(\cdot)$, $a_t \sim \pi(\cdot | s_t)$, $s_{t+1} \sim \mathbb{P}(\cdot | s_t, a_t)$, $r_{t+1} = r(s_{t+1} | s_t, a_t)$,
 102 and $c_{t+1} = c(s_{t+1} | s_t, a_t)$. The *state value function* of π is defined as $V_\pi(s) = \mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s]$.
 103 The goal of reinforcement learning is to maximize $J(\pi)$ defined as: $J(\pi) = \mathbb{E}_{s \sim \rho_0(\cdot)}[V_\pi(s)]$.
 104 The *cost-return* is: $J^c(\pi) = \mathbb{E}_{s \sim \rho_0(\cdot)}[\sum_{t=0}^{\infty} \gamma^t c_{t+1} | s_0 = s]$, we define the feasible policy set Π_C :
 105 $\Pi_C = \{\pi \in \Pi \text{ and } J^c(\pi) \leq b\}$. The goal of CMDP is to search the optimal policy π_* such that

$$\pi_* = \arg \max_{\pi \in \Pi_C} J(\pi). \quad (1)$$

106 **4 Safety Tasks**

107 In this section, we present four safety robot tasks, including robots with limited velocity, robots with
 108 circle, robots with goal and safety dexterousHands. The safety dexterousHands are new safe learning
 109 tasks, which are design to make deeper comparisons to the safe algorithms and measure research
 110 progress on safe RL. For more details about safety dexterousHands, see Appendix D.

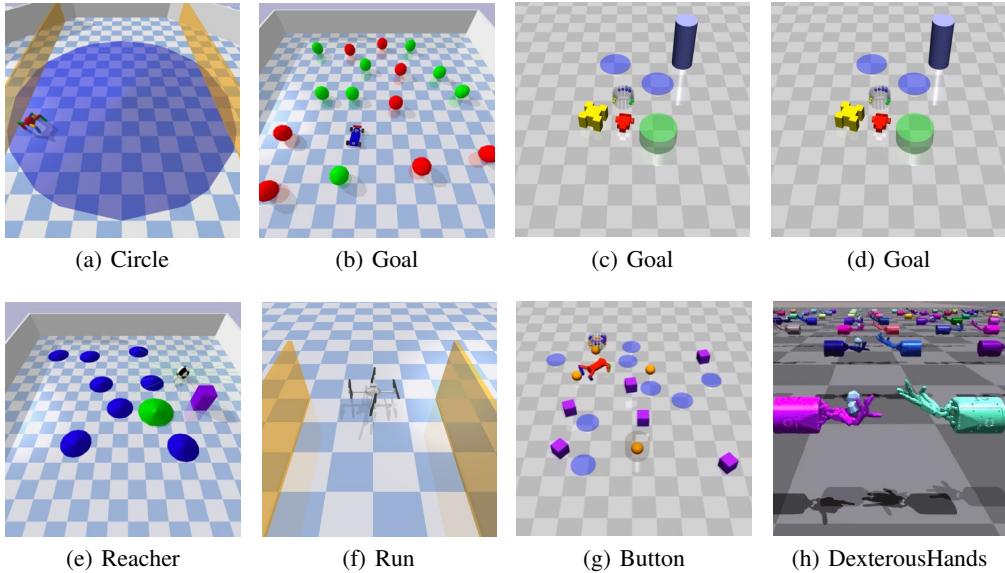


Figure 2: Fig.(a) Circle [1]: The robot is restricted to move in a given circle, but it is constrained to stay within a safe region more minor than the radius of the circle. Fig.(b), (c), (d) Goal [43]: Move the robot to a series of goal positions. When a goal is achieved, the goal location is randomly reset to someplace new while keeping the rest of the layout the same. Fig.(e) Reacher [27]: As soon the robot enters the goal zone, the goal is re-spawned such that the agent has to reach a series of goals. Fig.(f) Run [13]: Robots are rewarded for running through an avenue between two safety boundaries. Fig.(g) Button [43]: Press a series of goal buttons. Fig.(h) DexterousHands: Its task is to pass the ball from the right hand to the left hand, but some joints of the hand are freezing.

111 4.1 Task 1: Robots with Limited Velocity

112 We present the task of the robots with limited velocity according to [61, 60], where we train the
 113 robots to walk but imposed a limited velocity to the robots. We consider the tasks from MuJoCo
 114 [9]: Swimmer-v3, Walker2d-v3, Hopper-v3, HalfCheetah-v3, and Ant-v3. For the robots move on a
 115 two-dimensional plane, we calculate the cost as follows, $c(s, a) = \sqrt{v_x^2 + v_y^2}$; for the robots move
 116 move along a straight line, we calculate the cost as follows, $c(s, a) = |v_x|$, where v_x, v_y are the
 117 velocities of the agent in the x and y directions respectively.

118 4.2 Task 2: Robots with Circle

119 In the circle tasks, the goal is for an agent to move along the circumference of a circle while remaining
 120 within a safety region smaller than the radius of the circle. The reward and cost functions are:

$$R(s) = \frac{-yv_x + xv_y}{1 + |\sqrt{x^2 + y^2} - d|}, \quad C(s) = \mathbb{I}(|x| > x_{\lim}),$$

121 where x, y are the positions of the agent on the plane, v_x, v_y are the velocities of the agent along the
 122 x and y directions, d is the radius of the circle, and x_{\lim} specifies the range of the safety region. We
 123 set these parameters as follows: for Bullet Circle $d = 7$, $x_{\lim} = 4.5$; for Humanoid circle $d = 10$,
 124 $x_{\lim} = 2.5$, and for Ant circle $d = 10$, $x_{\lim} = 3$.

125 4.3 Task 3: Robots with Goal

126 In this benchmark, we consider four tasks of Goal: see Figure 2 (b)-(g), where those tasks share a
 127 common character that the robots move to a series of goal positions. We have presented some basic
 128 details of those tasks in the caption of Figure 2. Now, we present the task of Button according to
 129 [43]. Several immobile “buttons” are scattered throughout the environment, and the agent should
 130 navigate to and press the currently-highlighted Button. After the agent presses the correct Button, the
 131 environment will select and highlight a new goal button, keeping everything else fixed.

132 **4.4 Task 4: Safety DexterousHands**

133 A robot consists of two fixed hands, where the right hand holds a ball, its task is to pass the ball to
 134 the left hand. At the beginning, the ball falls randomly in the area around the right side, the right
 135 hand grabs the ball, and passes the object to the left hand. The base of the hand is fixed. In this task,
 136 the hand that holds the object initially cannot directly touch the target, nor can it directly roll the
 137 object to the other hand, so the object must be thrown up and stays in the air in the process. There are
 138 211-dimensional observations and 40-dimensional actions.

139 For each round t , let $x_{b,t}$ be the position of the ball and $x_{g,t}$ be the position of the goal, we use
 140 d_t to denote the distance between the ball and goal $d_t = \|x_{b,t} - x_{g,t}\|_2$. Let $d_{a,t}$ denote the
 141 angular position distance between the object and the goal, and the rotational difference $d_{r,t} =$
 142 $2 \arcsin \min\{|d_{a,t}|, 1.0\}$. The reward is defined as follows,

$$r_t = \exp\{-0.2(\alpha d_{d_t} + d_{r,t})\}, \quad (2)$$

143 where α is a constant balances translational and rotational rewards.

144 **Safety Joint.** See Figure 11, in this task, we constrain the freedom of the joint ④ of forefinger with a
 145 given region. The freedom angle of the joint ④ is among $[-15^\circ, 90^\circ]$, the safety task is to restrict to
 146 joint ④ within $[-7.5^\circ, 45^\circ]$. Let the ang_4 be the angle the joint ④ plays, and the cost is defined as:

$$c_t = \mathbb{I}(\text{ang_4} \notin [-7.5^\circ, 45^\circ]). \quad (3)$$

147 **Safety Finger.** In this task, we constrain the freedom of the joint ②, ③ and ④ of forefinger with a
 148 given region. The freedom of the joint ③ and ④ is among $[0^\circ, 90^\circ]$, the safety task is to restrict to
 149 joint ④ within $[0^\circ, 45^\circ]$. Let the ang_2 , ang_3 be the angle the joint ②, ③ plays, and the cost is:

$$c_t = \mathbb{I}(\text{ang_2} \notin [0^\circ, 45^\circ], \text{ or } \text{ang_3} \notin [0^\circ, 45^\circ], \text{ or } \text{ang_4} \notin [-7.5^\circ, 45^\circ]). \quad (4)$$

150 **5 Algorithms**

151 In this benchmark, we consider representative safe policy optimization algorithms. Appendix A
 152 contains all the key steps of those algorithms.

153 **CPO (Constrained Policy Optimization).** CPO [1] suggests to replace the reward objective $J(\pi)$
 154 and the cost constraint $J^c(\pi)$ with surrogate functions, which evaluates the objective $J(\pi)$ and
 155 constraint $J^c(\pi)$ according to the samples collected from the current policy π_k as follows,

$$\pi_{k+1} = \arg \max_{\pi \in \Pi_\theta} \mathbb{E}_{s \sim d_{\pi_k}^{\rho_0}(\cdot), a \sim \pi_\theta(\cdot|s)} [A_{\pi_k}(s, a)] \quad (5)$$

$$\text{s.t. } J^c(\pi_k) + \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\pi_k}^{\rho_0}(\cdot), a \sim \pi_\theta(\cdot|s)} [A_{\pi_k}^c(s, a)] \leq b, \quad (6)$$

$$\bar{D}_{\text{KL}}(\pi_\theta, \pi_k) = \mathbb{E}_{s \sim d_{\pi_k}^{\rho_0}(\cdot)} [\text{KL}(\pi_\theta, \pi_k)[s]] \leq \delta, \quad (7)$$

156 where π_k is short for π_{θ_k} . For the practical implementation, CPO introduces convex approximations
 157 to replace the term $A_{\pi_k}(s, a)$, $A_{\pi_k}^c(s, a)$, and $\bar{D}_{\text{KL}}(\pi_\theta, \pi_k)$, which is similar to TRPO [46].

158 **PCPO (Projection-based Constrained Policy Optimization).** PCPO [57] is an iterative method
 159 for optimizing policies in a two-step process: the first step performs a local reward improvement
 160 update, while the second step reconciles any constraint violation by projecting the policy back onto
 161 the constraint set. For the practical implementation, CPO inherits the key idea from TRPO [46] and
 162 CPO [1], where PCPO also use convex approximations to the reward objective and cost constraint.

163 **FOCOPS (First Order Constrained Optimization in Policy Space).** FOCOPS [61] finds the
 164 optimal update policy by solving a constrained optimization problem in the non-parameterized policy
 165 space, then it projects the update policy back into the parametric policy space.

166 **TRPO-L and PPO-L.** The Lagrangian approach is a standard way to solve CMDP (1), which is also
 167 known as primal-dual policy optimization:

$$(\pi_*, \lambda_*) = \arg \min_{\lambda \geq 0} \max_{\pi \in \Pi_\theta} \{J(\pi) - \lambda(J^c(\pi) - b)\}. \quad (8)$$

168 TRPO-Lagrangian and PPO-Lagrangian combine the Lagrangian approach with TRPO and PPO.
169 Concretely, PPO using the following clip term to replace $J(\pi)$ in (8),

$$\mathcal{L}_{\text{clip}}^r = \mathbb{E}_{s \sim d_{\pi_k}^{p_0}(\cdot), a \sim \pi_k(\cdot|s)} \left[-\min \left\{ \frac{\pi_{\theta}(a|s)}{\pi_k(a|s)} A_{\pi_k}(s, a), \text{clip} \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A_{\pi_k}(s, a) \right\} \right].$$

170 With $A_{\pi_k}(s, a)$ replacing $A_{\pi_k}^c(s, a)$ respectively, and obtain $\mathcal{L}_{\text{clip}}^c$. Let $\mathcal{L}_{\text{clip}}^c$ replace $J(\pi)$, PPO-L
171 updates the policy according to first-order optimizer. TRPO-L shares a similar idea but it is adaptive
172 to TRPO, due to the limitation of space, we provide the details in Appendix A.5.

173 **P3O (Penalized Proximal Policy Optimization).** P3O [60] solves the cumbersome constrained
174 policy iteration via a single minimization of an equivalent unconstrained problem:

$$\pi_{k+1} = \arg \min_{\pi \in \Pi_{\theta}} \left\{ \mathbb{E}_{s \sim d_{\pi_k}^{p_0}(\cdot), a \sim \pi_k(\cdot|s)} \left[\frac{\pi(a|s)}{\pi_k(a|s)} A_{\pi_k}(s, a) \right] + \kappa B(\pi, b) \right\}, \quad (9)$$

175 where κ is a positive scalar, and the penalty term $B(\pi, b)$ is defined as follows,

$$B(\pi, b) = \max \left\{ 0, \mathbb{E}_{s \sim d_{\pi_k}^{p_0}(\cdot), a \sim \pi_k(\cdot|s)} \left[\frac{\pi(a|s)}{\pi_k(a|s)} A_{\pi_k}^c(s, a) \right] + (1 - \gamma) (J^c(\pi_k) - b) \right\}. \quad (10)$$

176 P3O utilizes a simple yet effective penalty approach to eliminate cost constraints and removes the
177 trust-region constraint by the clipped surrogate objective.

178 6 Empirical Observations

179 In this section, we present the key factors influence the performance of the safety policy learning on a
180 comprehensive empirical comparison over 33 different safety robot tasks setting.

181 6.1 What Algorithm is Winner?

182 The first observation is that *there is no unambiguous winner*: no single algorithm performs best
183 among all the tasks for both reward performance and cost limit in all of the safety tasks, see Figure
184 3-6. In general, the approaches implemented by a deep neural network (e.g., P3O) achieve a better
185 performance than the convex approximation method (e.g., CPO); such an empirical observation
186 is sharper when the task is more complex (e.g., Safety DexterousHands, see Figure 6). For more
187 comprehensive empirical comparisons, see Appendix C.

188 CPO, TRPO-L, and PCPO obtain a relatively worse empirical performance than other algorithms.
189 This is not an accidental phenomenon due to the essential implementation mechanism of CPO, TRPO-
190 L, and PCPO. In practice, CPO and PCPO use convex approximations to replace $A_{\pi_{\theta_k}}(s, a)$ and
191 $\bar{D}_{\text{KL}}(\pi_{\theta}, \pi_{\theta_k})$ Eq.(5)-(7). Concretely, CPO suggests to use first-order Taylor expansion to replace (5)-
192 (6), uses second-order approximation to replace (7). Those convex approximations turn a non-convex
193 problem (5)-(7) to be a convex problem, it seems to make a simple solution, but this approach results
194 in many error sources and troubles in practice. It still lacks a theory analysis to show the difference
195 between the non-convex problem (5)-(7) and its convex approximations. Policy optimization is a
196 typical non-convex problem [56, 55]. Those convex approximations may introduce some errors for
197 the original issue. TRPO-L and PCPO inherit key ideas concerning convex approximations from
198 CPO, making TRPO-L and PCPO perform relatively worse in some complex safety tasks.

199 Instead of using a convex approximation to the objective function, FOCOPS, PPO, and P3O directly
200 optimize the surrogate objective function via the first-order method, and it does not depend on
201 any convex approximation. Additionally, the first-order method effectively avoids the expensive
202 computation for the inverse Fisher information matrix in CPO, PCPO, and TRPO-L. Those have
203 partially explained why FOCOPS, PPO, and P3O obtain better performance.

204 6.2 Comments on Robots with Limited Velocity

205 An important observation is that the safety task of robots with limited velocity is an *ill defined task*,
206 where a contradictory phenomenon occurs between reward improvement and constrained limited cost
207 for the robots like Ant-v3, Hopper-v3, HalfCheetah-v3, and Walker2d-v3.

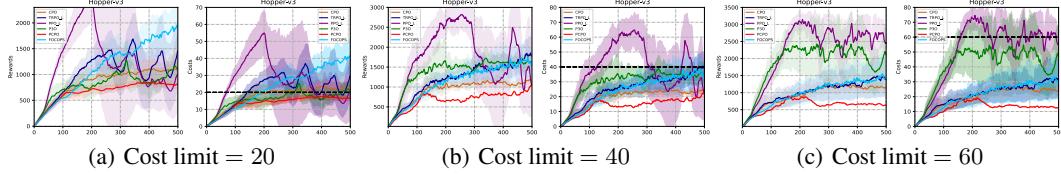


Figure 3: Learning curves for Hopper-v3 with limited velocity tasks. The shaded region represents the standard deviation of the score over three trials. Curves are smoothed uniformly for visual clarity. The x -axis represents the epoch, and each epoch updates the policy with 80×10^4 iterations.

208 Concretely, for example, the reward of Hopper-v3 is mainly determined by the following term ¹

$$r_t(s, a) \propto \frac{\|(x_{t'}, y_{t'}) - (x_t, y_t)\|_2}{t' - t} \approx \sqrt{v_x^2 + v_y^2} = c_t(s, a), \quad (11)$$

209 where (x_t, y_t) is the two-dimensional coordinate of the robot at time t , and $r_t(s, a)$ denotes the
210 reward with respect to the robot transforms from (x_t, y_t) to $(x_{t'}, y_{t'})$ after playing action a . The
211 relationship (11) shows that the reward function is proportional to the cost function. Recall the goal
212 of reinforcement learning is to search a policy improves the cumulative reward while controls the
213 boundedness of the cumulative cost with a given limit b , which is very difficult to be reconciled
214 between the reward improvement and the bounded cost limit. Additionally, from Eq.(11), we also
215 know that the learning curve of reward performance and cost performance can be very similar
216 potentially. The result of Figure 3 is consistent with our analysis, where the reward and cost learning
217 curve shares a very similar shape. Besides, from Figure 3, we know the reward increases while the
218 cost also increases, which is undesirable for safety learning.

219 Thus, the safety task of robots with limited velocity is ill-defined; for more empirical results, see
220 Figure 7. Such a worse case also occurs in the robots of Ant-v3, HalfCheetah-v3, and Walker2d-v3
221 since those robots also share a similar formulation (11). For more discussion, see Appendix B.

222 6.3 Comments on Lagrange Multiplier

223 We observe that the Lagrange parameter plays an important role for the agent in learning a safe policy.
224 If the policy violates the cost limit, the Lagrange parameter will be active and help the agent update
225 the policy alongs with the safe direction. Otherwise, the Lagrange parameter keeps calm, which
226 allows the agent improves the reward performance within the safety region.

227 Recall (8), the Lagrange method (e.g., PPO-L or TRPO-L) consider the following update rule of λ :

$$\lambda_{k+1} = \{\lambda_k + \eta(\hat{J}_k^c - b)\}_+, \quad (12)$$

228 where \hat{J}_k^c is an estimator for the cost function, $\eta > 0$ is step-size. Those empirical result shown
229 in Figure 4 is consistent with the update rule of λ in Eq.(12). If the estimated cost function is
230 under the target threshold b , then λ keeps calm, which implies λ is not activated. Such an empirical
231 phenomenon gives significant expression to the Humanoid environment. While if the estimated cost
232 exceeds the target threshold b , λ will be activated, which requires the agent to play a policy on the
233 safe region. This implies Lagrange parameter plays an important role for the agent to learn a safe
234 policy. Additionally, from the above discussion, we know Figure 4 provides a visualization way to
235 show the difficulty of different tasks, where the task actives much quantification of λ , such a task is
236 more challenging to obtain a safe policy since the agent needs more trials and errors.

237 6.4 Sparse Cost

238 We observe that it is very challenging for the agent to learn a safe policy if the cost is too sparse.
239 From Figure 5, in the SafetyDroneCircle, TRPO-L, FOCOPS, CPO, and PCPO achieve a very low
240 reward performance. This is because the total cost limit is tiny, so the robot obtains a small cost for
241 each step, implying the robot knows very little about the environment. Thus the robot pays much
242 effort into exploring the environment to find a safe policy.

¹We present the details with respect to reward according to the open implementation: <https://www.gymlibrary.ml/environments/mujoco/hopper/>.

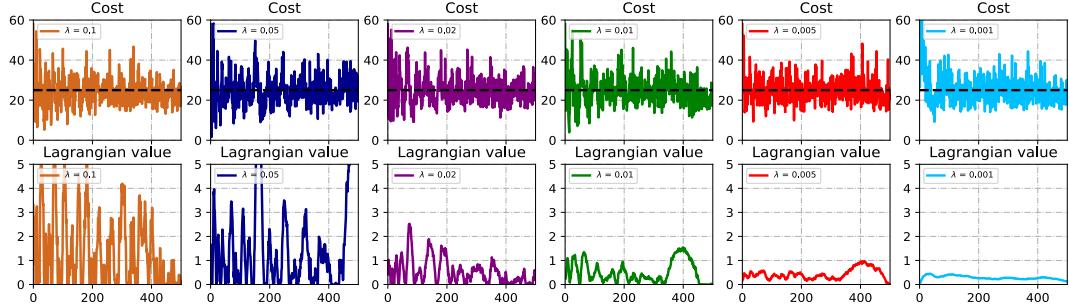


Figure 4: Cost constraint with respect to the Lagrangian hyper-parameter λ .

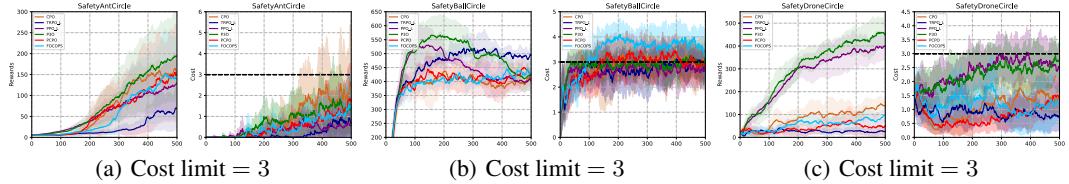


Figure 5: Learning curves for SafetyAntCircle, SafetyBallCircle, and SafetyDroneCircle. The shaded region represents the standard deviation of the score over three trials. Curves are smoothed uniformly for visual clarity.

243 A more general observation is that the cost helps the robot search for a safe policy. From Figure 5,
 244 we know the algorithm near the cost limit usually obtains better performance than that far away from
 245 the cost limit. This is because the higher cost implies the robot learns more about the environment,
 246 which helps the robot understand a safer policy. This observation further illustrates that cost plays a
 247 quantization measurement of exploration. Additionally, we know that if the cost limit is too large, the
 248 task could be reduced to an unconstrained policy learning since the robot explores the environment
 249 with a huge cost. Thus, we should consider the trade-off between a small and large cost limit for
 250 real-world applications to design a safety policy learning algorithm.

251 7 Safe Robot Learning in Dexterous Hands

252 Dexterous manipulation with hands to complete complex task is one of the most challenging problem
 253 in the robotics [38]. In this section, we propose the Safety DexterousHands tasks, which is a new and
 254 challenging robot task for measuring research progress on safe reinforcement learning. For a deeper
 255 comparison to the safe algorithms, in this section we present the performance with two settings:
 256 single-agent and multi-agent.

257 In contrast to single agent RL , in Multi-Agent Reinforcement Learning (MARL) [28], it's more
 258 difficult to optimise policy because of the instability of multi-agent systems. In MARL settings, when
 259 an agent optimises its policy, it needs to consider other agents' policy simultaneously. When it comes
 260 to safe MARL settings, it's more challenging to improve policy performance than purely reward
 261 optimisation, since it not only needs to improve agents' reward, but also needs to guarantee agents'
 262 safety. We provide a safe MARL baseline benchmark for safe MARL research on challenging tasks
 263 of safety DexterousHands, in which the MACPO [28], MAPPO-L [28], MAPPO [58], HAPPO [32],
 264 IPPO [17] are all implemented to investigate the safety and reward performance. MACPO is developed
 265 based CPO and HATRPO [32] by leveraging contrained trust region optimisation and multi-agent
 266 sequential decision theory, MACPO can achieve hard constrained optimisation to ensure safety for
 267 multi-agent systems; MAPPO-L is developed based on primal-dual methods, which can achieve soft
 268 constrained policy optimisation to improve the safety of multi-agent systems.

269 **Single-Agent Safety Learning.** Figure 6 and Table 4 show the results of single-agent safe RL
 270 algorithm over the Safety Joint and Safety Finger tasks. In general, although the tasks are all
 271 two-hand cooperation, Safety Joint outperforms Safety Finger in most cases, which implies finger
 272 manipulation is more difficult than joint manipulation. This observation seems to oppose the human
 273 behavior: a person makes a better finger manipulation than a single joint, which implies that a task

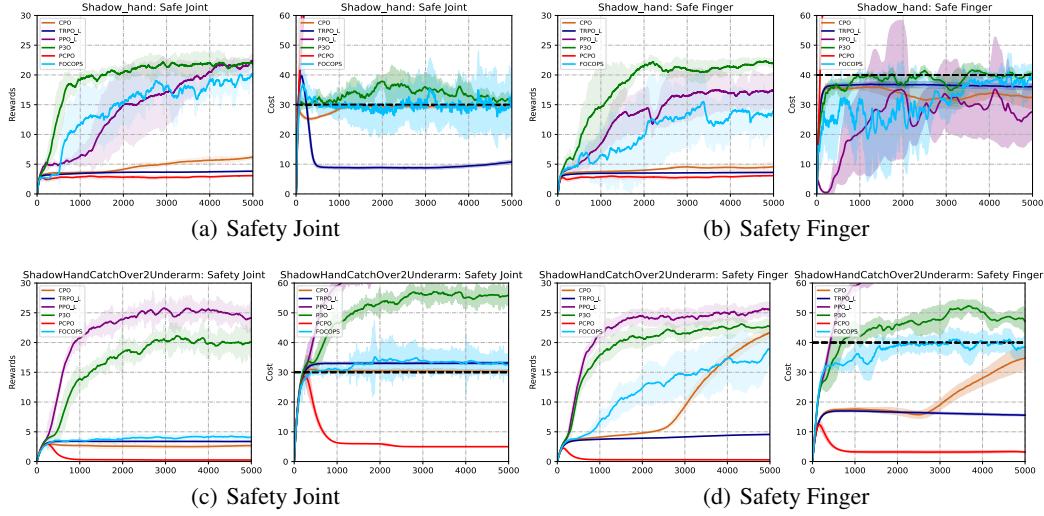


Figure 6: Learning curves for Single-agent Safety DexterousHands. The shaded region represents the standard deviation of the score over three trials. Curves are smoothed uniformly for visual clarity.

is difficult for human but it is simple for a machine. Besides, we observe that FOCOPS and P3O outperforms CPO and PCPO significantly, this is because the previous task is relatively easy. The advantage of our VRMPO becomes more powerful when the task is more difficult. CPO and PCPO obtain a relatively worse empirical performance than other algorithms. This is not an accidental phenomenon due to the essential implementation mechanism of CPO and PCPO. Instead of using a convex approximation for the objective function, FOCOPS and P3O directly optimizes the surrogate objective function via the first-order method, and it does not depend on any convex approximation.

Multi-Agent Safety Learning. The empirical study reveals the findings that MACPO and MAPPO-lagrangian algorithms can achieve comparable reward performance compared with the state-of-art MARL algorithms while guaranteeing safety of agents. Figure 12 shows performance comparisons on tasks of Safe ShadowHandOver 2x6 task two in terms of cost and reward, the results indicate that MACPO and MAPPO-lagrangian algorithms are the only ones that learn the safety constraints compared with MAPPO, HAPPO and IPPO, while achieving satisfying performance in terms of the reward. This is consistent with what has been found in [28] work.

8 Conclusion and Future Work

In this study, we standardize the safe policy optimization method for solving constrained Markov decision processes and introduced the first benchmark named SafePO. We implemented constrained policy optimization (CPO) type algorithms and Lagrangian based algorithms and test them on four kinds of safety-aware robotic tasks. Moreover, we design Safety DexterousHands tasks for safe multi-agent RL tasks. Importantly, we design SafePO in a way that it can collects samples in parallel and synchronizes the policy gradients on different physical CPU cores, this makes SafePO highly efficient for training a safe policy. On top of the benchmark comparison, we also highlight the interplay of different attributes on the safety performance, and illustrate the difficulty of safety learning on the sparse cost.

In this benchmark, we only consider the safety policy learning under the framework of constrained Markov decision process, which is one category of safe reinforcement learning [25]. In the future, we will consider the following problems for safety policy learning. 1) **Zero Constraint Violation:** In real-world applications, for example, autonomous vehicle [53, 36] or power systems [26, 52] it is catastrophic if the system plays violations of constraints [5]. Thus, achieving its goal guarantees zero constraint violation, an important problem for safety learning. 2) **State-wise Safety Learning.** This benchmark only considers relaxed trajectory-wise constraints. However, those approaches may fail to satisfy the hard state-wise safety constraints [19], where the state-wise constraint satisfaction is more practical than safe trajectories, which yields a constrained optimization problem at each decision-making step.

308 **References**

- 309 [1] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization.
310 In *Proceedings of International Conference on Machine Learning (ICML)*, volume 70, pages
311 22–31, 2017.
- 312 [2] Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and
313 Ufuk Topcu. Safe reinforcement learning via shielding. In *Proceedings of the AAAI Conference
314 on Artificial Intelligence*, volume 32, 2018.
- 315 [3] Eitan Altman. *Constrained Markov decision processes*. CRC Press, 1999.
- 316 [4] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
317 Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 318 [5] Qinbo Bai, Amrit Singh Bedi, Mridul Agarwal, Alec Koppel, and Vaneet Aggarwal. Achieving
319 zero constraint violation for constrained reinforcement learning via primal-dual approach. *arXiv
320 preprint arXiv:2109.06332*, 2021.
- 321 [6] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich
322 Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, Julian Schrittwieser,
323 Keith Anderson, Sarah York, Max Cant, Adam Cain, Adrian Bolton, Stephen Gaffney, He-
324 len King, Demis Hassabis, Legg Shane, and Stig Petersen. Deepmind lab. *arXiv preprint
325 arXiv:1612.03801*, 2016.
- 326 [7] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning
327 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence
328 Research*, 47:253–279, 2013.
- 329 [8] Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. Safe model-based
330 reinforcement learning with stability guarantees. *Advances in neural information processing
331 systems*, 30, 2017.
- 332 [9] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,
333 and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 334 [10] Lukas Brunke, Melissa Greeff, Adam W Hall, Zhaocong Yuan, Siqi Zhou, Jacopo Panerati, and
335 Angela P Schoellig. Safe learning in robotics: From learning-based control to safe reinforcement
336 learning. *Annual Review of Control, Robotics, and Autonomous Systems*, 5, 2021.
- 337 [11] Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained
338 reinforcement learning with percentile risk criteria. *The Journal of Machine Learning Research*,
339 18(1):6070–6120, 2017.
- 340 [12] Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. A
341 lyapunov-based approach to safe reinforcement learning. In *Advances in Neural Information
342 Processing Systems (NeurIPS)*, 2018.
- 343 [13] Yinlam Chow, Ofir Nachum, Aleksandra Faust, Edgar Duenez-Guzman, and Mohammad
344 Ghavamzadeh. Lyapunov-based safe policy optimization for continuous control. *arXiv preprint
345 arXiv:1901.10031*, 2019.
- 346 [14] Yinlam Chow, Ofir Nachum, and Mohammad Ghavamzadeh. Path consistency learning in tsallis
347 entropy regularized mdps. In *International Conference on Machine Learning (ICML)*, pages
348 978–987, 2018.
- 349 [15] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying
350 generalization in reinforcement learning. In *International Conference on Machine Learning*,
351 pages 1282–1289. PMLR, 2019.
- 352 [16] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games,
353 robotics and machine learning (2016–2019), 2019.

- 354 [17] Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS
 355 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft
 356 multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- 357 [18] Brian Delhaisse, Leonel Rozo, and Darwin G Caldwell. Pyrobolearn: A python framework for
 358 robot learning practitioners. In *Conference on Robot Learning*, pages 1348–1358. PMLR, 2020.
- 359 [19] Priya L Donti, David Rolnick, and J Zico Kolter. Dc3: A learning method for optimization with
 360 hard constraints. *arXiv preprint arXiv:2104.12225*, 2021.
- 361 [20] Yan Duan, Xi Chen, Rein Houthooft, John Schulman, and Pieter Abbeel. Benchmarking
 362 deep reinforcement learning for continuous control. In *International Conference on Machine
 363 Learning (ICML)*, pages 1329–1338, 2016.
- 364 [21] Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. Challenges of real-world reinforce-
 365 ment learning. *arXiv preprint arXiv:1904.12901*, 2019.
- 366 [22] Zackory Erickson, Vamsee Gangaram, Ariel Kapusta, C Karen Liu, and Charles C Kemp. As-
 367 ssistive gym: A physics simulation framework for assistive robotics. In *2020 IEEE International
 368 Conference on Robotics and Automation (ICRA)*, pages 10169–10176. IEEE, 2020.
- 369 [23] Linxi Fan, Yuke Zhu, Jiren Zhu, Zihua Liu, Orien Zeng, Anchit Gupta, Joan Creus-Costa, Silvio
 370 Savarese, and Li Fei-Fei. Surreal: Open-source reinforcement learning framework and robot
 371 manipulation benchmark. In *Conference on Robot Learning*, pages 767–782. PMLR, 2018.
- 372 [24] C Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier
 373 Bachem. Brax—a differentiable physics engine for large scale rigid body simulation. *arXiv
 374 preprint arXiv:2106.13281*, 2021.
- 375 [25] Javier García and Fernando Fernández. A comprehensive survey on safe reinforcement learning.
 376 *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- 377 [26] Istemihan Genc, Ruisheng Diao, Vijay Vittal, Sharma Kolluri, and Sujit Mandal. Decision
 378 tree-based preventive and corrective control applications for dynamic security enhancement in
 379 power systems. *IEEE Transactions on Power Systems*, 25(3):1611–1619, 2010.
- 380 [27] Sven Gronauer. Bullet-safety-gym: A framework for constrained reinforcement learning. 2022.
- 381 [28] Shangding Gu, Jakub Grudzien Kuba, Munning Wen, Ruiqing Chen, Ziyan Wang, Zheng Tian,
 382 Jun Wang, Alois Knoll, and Yaodong Yang. Multi-agent constrained policy optimisation. *arXiv
 383 preprint arXiv:2110.02793*, 2021.
- 384 [29] Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang,
 385 and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications.
 386 *arXiv preprint arXiv:2205.10330*, 2022.
- 387 [30] Jemin Hwangbo, Joonho Lee, Alexey Dosovitskiy, Dario Bellicoso, Vassilios Tsounis, Vladlen
 388 Koltun, and Marco Hutter. Learning agile and dynamic motor skills for legged robots. *Science
 389 Robotics*, 4(26):eaau5872, 2019.
- 390 [31] Torsten Koller, Felix Berkenkamp, Matteo Turchetta, and Andreas Krause. Learning-based
 391 model predictive control for safe exploration. In *Conference on Decision and Control (CDC)*,
 392 pages 6059–6066. IEEE, 2018.
- 393 [32] Jakub Grudzien Kuba, Ruiqing Chen, Munning Wen, Ying Wen, Fanglei Sun, Jun Wang, and
 394 Yaodong Yang. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv
 395 preprint arXiv:2109.11251*, 2021.
- 396 [33] Hoang Le, Cameron Voloshin, and Yisong Yue. Batch policy learning under constraints. In
 397 *International Conference on Machine Learning (ICML)*, pages 3703–3712, 2019.
- 398 [34] Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning
 399 quadrupedal locomotion over challenging terrain. *Science robotics*, 5(47):eabc5986, 2020.

- 400 [35] Jan Leike, Miljan Martic, Victoria Krakovna, Pedro A Ortega, Tom Everitt, Andrew Lefrancq,
401 Laurent Orseau, and Shane Legg. Ai safety gridworlds. *arXiv preprint arXiv:1711.09883*, 2017.
- 402 [36] Quanyi Li, Zhenghao Peng, Zhenghai Xue, Qihang Zhang, and Bolei Zhou. Metadrive: Com-
403 posing diverse driving scenarios for generalizable reinforcement learning. *arXiv preprint
arXiv:2109.12674*, 2021.
- 405 [37] Zuxin Liu, Zijian Guo, Zhepeng Cen, Huan Zhang, Jie Tan, Bo Li, and Ding Zhao. On the
406 robustness of safe reinforcement learning under observational perturbations. *arXiv preprint
arXiv:2205.14691*, 2022.
- 408 [38] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
409 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High
410 performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*,
411 2021.
- 412 [39] Takahiro Miki, Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco
413 Hutter. Learning robust perceptive locomotion for quadrupedal robots in the wild. *Science
Robotics*, 7(62):eabk2822, 2022.
- 415 [40] Teodor Mihai Moldovan and Pieter Abbeel. Safe exploration in markov decision processes.
416 *arXiv preprint arXiv:1205.4810*, 2012.
- 417 [41] Santiago Paternain, Luiz FO Chamon, Miguel Calvo-Fullana, and Alejandro Ribeiro. Con-
418 strained reinforcement learning has zero duality gap. In *Advances in Neural Information
Processing Systems (NeurIPS)*, 2019.
- 420 [42] Martin Pecka and Tomas Svoboda. Safe exploration techniques for reinforcement learning—an
421 overview. In *International Workshop on Modelling and Simulation for Autonomous Systems*,
422 pages 357–375. Springer, 2014.
- 423 [43] Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking Safe Exploration in Deep
424 Reinforcement Learning. 2019.
- 425 [44] Reazul Hasan Russel, Mouhacine Benosman, and Jeroen Van Baar. Robust constrained-
426 mdps: Soft-constrained robust policy optimization under model uncertainty. *arXiv preprint
arXiv:2010.04870*, 2020.
- 428 [45] Harsh Satija, Philip Amortila, and Joelle Pineau. Constrained markov decision processes via
429 backward value functions. In *International Conference on Machine Learning (ICML)*, pages
430 8502–8511, 2020.
- 431 [46] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust
432 region policy optimization. In *International Conference on Machine Learning (ICML)*, pages
433 1889–1897, 2015.
- 434 [47] Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning
435 by pid lagrangian methods. In *International Conference on Machine Learning*, pages 9133–9143.
436 PMLR, 2020.
- 437 [48] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David
438 Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy P. Lillicrap, and Martin A.
439 Riedmiller. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- 440 [49] Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization.
441 *International Conference on Learning Representation (ICLR)*, 2019.
- 442 [50] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based
443 control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages
444 5026–5033. IEEE, 2012.
- 445 [51] Cameron Voloshin, Hoang M Le, Nan Jiang, and Yisong Yue. Empirical study of off-policy
446 policy evaluation for reinforcement learning. *arXiv preprint arXiv:1911.06854*, 2019.

- 447 [52] Thanh Long Vu, Sayak Mukherjee, Tim Yin, Renke Huang, Jie Tan, and Qiuhua Huang. Safe
448 reinforcement learning for emergency load shedding of power systems. In *2021 IEEE Power &*
449 *Energy Society General Meeting (PESGM)*, pages 1–5. IEEE, 2021.
- 450 [53] Lu Wen, Jingliang Duan, Shengbo Eben Li, Shaobing Xu, and Huei Peng. Safe reinforcement
451 learning for autonomous vehicles through parallel constrained policy optimization. In *2020*
452 *IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pages 1–7.
453 IEEE, 2020.
- 454 [54] Long Yang, Jiaming Ji, Juntao Dai, Yu Zhang, Pengfei Li, and Gang Pan. Cup: A conservative
455 update policy algorithm for safe reinforcement learning. *arXiv preprint arXiv:2202.07565*,
456 2022.
- 457 [55] Long Yang, Gang Zheng, Haotian Zhang, Yu Zhang, Qian Zheng, and Gang Pan. Policy
458 optimization with stochastic mirror descent. *Association for the Advancement of Artificial*
459 *Intelligence (AAAI)*, 2022.
- 460 [56] Long Yang, Qian Zheng, and Gang Pan. Sample complexity of policy gradient finding second-
461 order stationary points. In *Association for the Advancement of Artificial Intelligence (AAAI)*,
462 2021.
- 463 [57] Tsung-Yen Yang, Justinian Rosca, Karthik Narasimhan, and Peter J Ramadge. Projection-based
464 constrained policy optimization. In *International Conference on Learning Representation*
465 (*ICLR*), 2020.
- 466 [58] Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising
467 effectiveness of ppo in cooperative, multi-agent games. *arXiv preprint arXiv:2103.01955*, 2021.
- 468 [59] Zhaocong Yuan, Adam W Hall, Siqi Zhou, Lukas Brunke, Melissa Greeff, Jacopo Panerati,
469 and Angela P Schoellig. safe-control-gym: a unified benchmark suite for safe learning-based
470 control and reinforcement learning. *arXiv preprint arXiv:2109.06325*, 2021.
- 471 [60] Linrui Zhang, Li Shen, Long Yang, Shixiang Chen, Bo Yuan, Xueqian Wang, Dacheng Tao,
472 et al. Penalized proximal policy optimization for safe reinforcement learning. *Proceedings of*
473 *the International Joint Conference on Artificial Intelligence (IJCAI)*, 2022.
- 474 [61] Yiming Zhang, Quan Vuong, and Keith Ross. First order constrained optimization in policy
475 space. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, 2020.

476 **Checklist**

- 477 1. For all authors...
 - 478 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
479 contributions and scope? **[Yes]**
 - 480 (b) Did you describe the limitations of your work? **[Yes]**
 - 481 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
 - 482 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
483 them? **[Yes]**
- 484 2. If you are including theoretical results...
 - 485 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
 - 486 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 487 3. If you ran experiments...
 - 488 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
489 mental results (either in the supplemental material or as a URL)? **[Yes]**
 - 490 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
491 were chosen)? **[Yes]**
 - 492 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
493 ments multiple times)? **[Yes]**
 - 494 (d) Did you include the total amount of compute and the type of resources used (e.g., type
495 of GPUs, internal cluster, or cloud provider)? **[Yes]**
- 496 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 497 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
 - 498 (b) Did you mention the license of the assets? **[Yes]**
 - 499 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
 - 500 (d) Did you discuss whether and how consent was obtained from people whose data you're
501 using/curating? **[Yes]**
 - 502 (e) Did you discuss whether the data you are using/curating contains personally identifiable
503 information or offensive content? **[N/A]**
- 504 5. If you used crowdsourcing or conducted research with human subjects...
 - 505 (a) Did you include the full text of instructions given to participants and screenshots, if
506 applicable? **[N/A]**
 - 507 (b) Did you describe any potential participant risks, with links to Institutional Review
508 Board (IRB) approvals, if applicable? **[N/A]**
 - 509 (c) Did you include the estimated hourly wage paid to participants and the total amount
510 spent on participant compensation? **[N/A]**