# Feasible Policy Iteration for Safe Reinforcement Learning

Yujie Yang, Zhilong Zheng, Shengbo Eben Li, Wei Xu, Jingjing Liu, Xianyuan Zhan, Ya-Qin Zhang

*Abstract*—**Safety is the priority concern when applying reinforcement learning (RL) algorithms to real-world control problems. While policy iteration provides a fundamental algorithm for standard RL, an analogous theoretical algorithm for safe RL remains absent. In this paper, we propose feasible policy iteration (FPI), the first foundational dynamic programming algorithm for safe RL. FPI alternates between policy evaluation, region identification and policy improvement. This follows actor-critic-scenery (ACS) framework where scenery refers to a feasibility function that represents a feasible region. A region-wise update rule is developed for the policy improvement step, which maximizes state-value function inside the feasible region and minimizes feasibility function outside it. With this update rule, FPI guarantees monotonic expansion of feasible region, monotonic improvement of state-value function, and geometric convergence to the optimal safe policy. Experimental results demonstrate that FPI achieves strictly zero constraint violation on low-dimensional tasks and outperforms existing methods in constraint adherence and reward performance on high-dimensional tasks.**

*Index Terms*—**safe reinforcement learning, feasible policy iteration, actor-critic-scenery, monotonic improvement, geometric convergence.**

## I. INTRODUCTION

**R**EINFORCEMENT learning has achieved promising performance on many challenging tasks such as video games [1], board games [2], robotics [3], and autonomous driving [4]. RL solves an optimal control problem (OCP) by finding a policy that maximizes the expected cumulative rewards. However, many real-world control tasks require not only maximizing rewards but also satisfying safety constraints. Such problems can be formulated as constrained OCPs [5]. In this paper, we consider stepwise deterministic constraints, which require strict constraint satisfaction at every step.

Policy iteration (PI) is a fundamental dynamic programming algorithm for unconstrained RL [6]. It iteratively performs

Yujie Yang and Zhilong Zheng contributed equally to this work. All correspondence should be sent to Shengbo Eben Li.

Yujie Yang, Zhilong Zheng, and Shengbo Eben Li are with the School of Vehicle and Mobility, Tsinghua University, Beijing, 10084, China (e-mail: yangyj21@mails.tsinghua.edu.cn, zhengzl22@mails.tsinghua.edu.cn, lishbo@tsinghua.edu.cn).

Wei Xu is with the College of Artificial Intelligence, Tsinghua University, Beijing, 10084, China (e-mail: weixu@tsinghua.edu.cn).

Jingjing Liu, Xianyuan Zhan, and Ya-Qin Zhang are with the Institute for AI Industry Research, Tsinghua University, Beijing, 10084, China (e-mail: jjliu@air.tsinghua.edu.cn, zhanxianyuan@air.tsinghua.edu.cn, zhangyaqin@air.tsinghua.edu.cn).

two steps: policy evaluation (PEV), which computes the state-value function of the current policy, and policy improvement (PIM), which updates the policy by maximizing the state-value function. PI lays the theoretical foundation for many modern RL algorithms [7]–[9]. Despite its success, PI is limited to unconstrained RL and cannot handle safety constraints. In the context of safe RL, a foundational algorithm analogous to PI remains undeveloped.

Existing safe RL methods essentially combine constrained optimization techniques with RL algorithms. A prominent class of methods is called *iterative unconstrained RL*, which reformulates the safe RL problem as a sequence of unconstrained optimization problems and solves each one with standard RL algorithms. The most widely used reformulation method is the Lagrange multiplier method [10], where constraints are added to the objective function with Lagrange multipliers to form an unconstrained optimization problem. For specific types of constraints, such as the cost value function in the constrained Markov decision process (CMDP), the Lagrange function can be rearranged into a discounted summation. With this rearrangement, the unconstrained problem turns out to be an RL problem itself and can be solved using any RL algorithm. The problem with iterative unconstrained RL methods is that they suffer from slow convergence and unstable training. The slow convergence arises from the need to solve an RL problem in every iteration of dual ascent, which, in theory, makes these algorithms converge a magnitude slower than standard RL algorithms. The instability in training stems from the inherent characteristic of the Lagrange method, which only guarantees convergence at the end but provides no assurance for intermediate objective functions or constraints. In safe RL, this lack of intermediate guarantee is manifested as oscillations of reward performance and constraint violations during training.

Another class of methods is called *constrained policy optimization*, which follows the policy optimization framework of standard RL and incorporates safety constraints into each iteration. Policy optimization methods treat RL as an optimization problem and solve it using established optimization algorithms. For example, vanilla policy gradient uses a gradient descent method, and trust region policy optimization (TRPO) [11] uses a trust region method. These optimization algorithms typically involve an iteration process. In safe RL, the constraint can be incorporated into these iterations. The most representative example is constrained policy optimization (CPO) [12], which adapts TRPO's trust region update
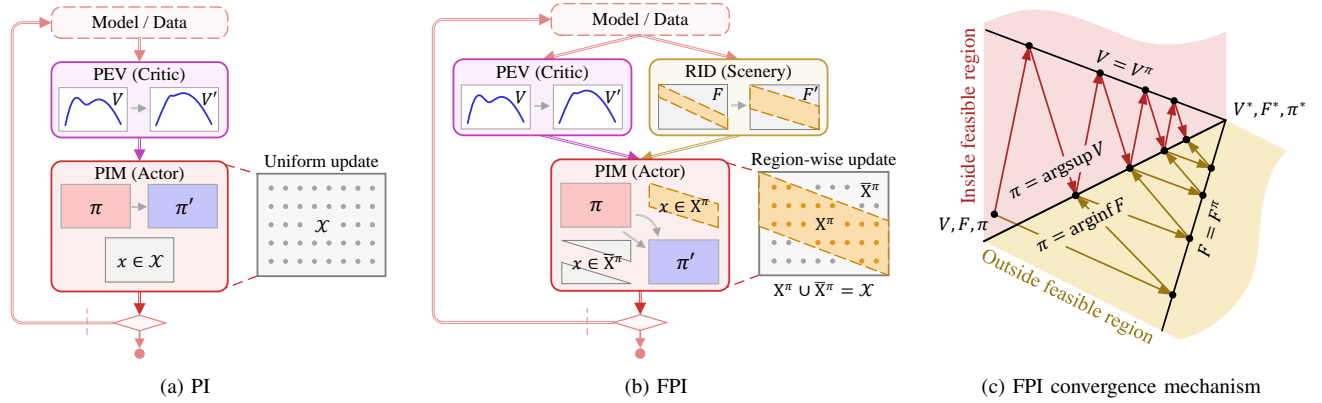
Fig. 1. Iteration framework and convergence mechanism of FPI. (a) PI adopts the AC framework. (b) The ACS framework for FPI features an additional RID step and a region-wise PIM. (c) The iteration of FPI converges to the optimal solution of safe RL.
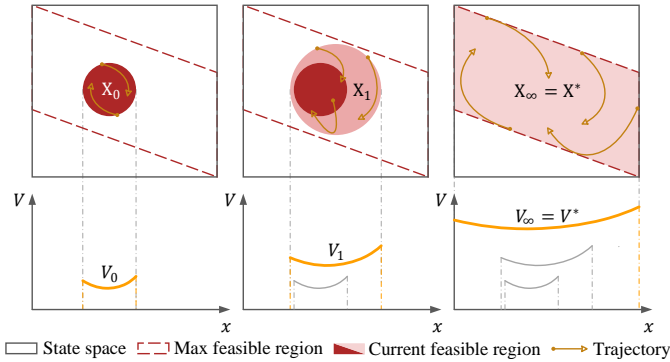


Fig. 2. Monotonicity and convergence of FPI. The feasible region monotonically expands and converges to the maximum feasible region. The state-value function monotonically increases inside the feasible region and converges to the optimal state-value function.

framework by adding a safety constraint in each iteration. The constraint function is linearized within the trust region to obtain an analytical solution. The problem with these methods is that they suffer from the infeasibility issue, i.e., they may fail to find a constraint-satisfying solution in some iterations. This is because the constraint, which remains fixed throughout training, is too stringent, especially at an early stage. At the beginning of training, the policy is usually randomly initialized and requires numerous updates before it can satisfy the safety constraint. When infeasibility occurs, these methods resort to alternative update rules, e.g., solely minimizing the constraint function [12]. This prevents the algorithm from optimizing rewards until the constraints are fully met, resulting in overly conservative behavior and low training efficiency.

The aforementioned challenges of existing safe RL algorithms stem from the lack of a theoretically effective and efficient iteration approach for safe RL that is analogous to PI for standard RL. To this end, we propose the first foundational dynamic programming algorithm for safe RL called feasible policy iteration (FPI). Different from PI which adopts an actor-critic (AC) framework, FPI follows an actor-critic-scenery (ACS) framework, which additionally includes a region identification (RID) step for updating the "scenery" and performs a region-wise update in PIM, as shown in Fig. 1. The

"scenery" in ACS refers to a feasibility function that represents the long-term constraint-satisfying set of a policy. Region-wise PIM updates the policy by minimizing the feasibility function outside the feasible region and maximizing the state-value function inside the feasible region. By recursively applying the self-consistency conditions, we prove that FPI guarantees monotonic expansion of the feasible region and monotonic increase of the state-value function within the feasible region. Furthermore, by deriving the contraction properties of two Bellman operators, we prove that both the feasibility function and state-value function geometrically converge to their optimum, as illustrated in Fig. 2. Experiment results show that FPI can learn strictly safe and near-optimal policies on low-dimensional tasks, and it can also achieve fewer violations and better performance on high-dimensional tasks.

## II. RELATED WORK

### A. Iterative unconstrained RL methods

Most iterative unconstrained RL methods use the Lagrange method and solve the dual problem using dual ascent, where the minimization step involves an unconstrained RL problem. The theoretical basis for these methods is established by Paternain et al. [13], who prove that the safe RL problem has zero duality gap. Under this framework, different types of constraints can be used. For example, Chow et al. [14] constrain the conditional value-at-risk of the cumulative cost, forming a probabilistic constraint. Tessler et al. [15] add the cost signal to the reward and treat the combined discounted sum as a new value function. Yu et al. [16] learn a Hamilton-Jacobi (HJ) reachability function and constrain the policy in its reachable set. Several works have also employed HJ reachability functions as safety constraints [17]–[20]. Besides HJ reachability function, other safety certificates frequently used in safe RL include control barrier function (CBF) [21]–[26], safety index [27], [28], and Lyapunov function [29]–[31]. For example, Yang et al. [24] learn a neural CBF and construct a multi-step barrier constraint for policy update. Ma et al. [28] learn a safety index with a given functional form and adjustable parameters by minimizing energy increase. Chow et al. [30] construct a Lyapunov function to guarantee the

global safety of the policy via a set of local linear constraints. The limitation of iterative unconstrained RL methods is that their convergence is slow because they introduce an additional iterative loop on top of standard RL. Moreover, their learning process exhibit significant oscillation [32] because of the lack of intermediate guarantees of the Lagrange method. While empirical techniques such as PID-Lagrangian [32], [33] are proposed to reduce oscillations, they still cannot provide theoretical guarantees for intermediate policies.

### B. Constrained policy optimization methods

The most representative algorithm of this class is CPO [12], which builds on ideas of trust region update and Taylor expansion approximation from TRPO [11]. Within this framework, CPO further approximates the constraint of the cost value function with an affine function and analytically solves the resulting constrained optimization problem. The authors prove that CPO ensures a lower bound on policy performance and an upper bound on constraint violation in every iteration. To avoid the computationally expensive line search in CPO, Yang et al. [34] propose projection-based policy optimization (PCPO), which first performs a reward improvement update and then projects the policy back onto the constrained set. Zhang et al. [35] propose first-order constrained optimization in policy space (FOCOPS), which first solves for the optimal update policy in the non-parameterized policy space and then projects it back into the parametric policy space. Following the policy projection route, Yang et al. [36] develop constrained update projection (CUP), which generalizes the surrogate functions estimator and only requires first-order optimizers without relying on any strong approximation on the convexity of the objectives. Inspired by the interior-point method, Liu et al. [37] incorporate a logarithmic barrier function of the constraint to the objective function and optimize it using a first-order policy optimization method. Additionally, several algorithms integrate the Lagrange multiplier method with policy optimization, such as proximal policy optimization with Lagrangian (PPO-Lagrangian) [38], TRPO-Lagrangian [38], and augmented PPO [39], also falling under this category. A limitation of these methods is their strict requirement for the policy to satisfy the original safety constraint in all iterations. Enforcing such constraints can lead to infeasibility issues, as mentioned by Achiam et al. [12]. It also causes overly conservative behavior in an early training stage, as the policy must focus solely on minimizing constraint violations without improving performance.

## III. PROBLEM FORMULATION

### A. Markov decision process with state constraints

In this paper, we consider a deterministic Markov decision process (MDP) described by a 5-tuple $(\mathcal{X}, \mathcal{U}, f, r, \gamma)$, where $\mathcal{X} \subseteq \mathbb{R}^n$ is the finite state space, $\mathcal{U} \subseteq \mathbb{R}^m$ is the finite action space, $f : \mathcal{X} \times \mathcal{U} \to \mathcal{X}$ is the dynamics model, $r : \mathcal{X} \times \mathcal{U} \to \mathbb{R}$ is the reward function, $0 < \gamma < 1$ is the discount factor. The policy in the MDP considered in this paper is a deterministic

function that maps a state to an action: $\pi : \mathcal{X} \to \mathcal{U}$. The state-value function under a policy $\pi$, $V^\pi : \mathcal{X} \to \mathbb{R}$, is defined as the discounted sum of rewards under $\pi$:

$$V^\pi(x_0) = \sum_{t=0}^{\infty} \gamma^t r(x_t, \pi(x_t)), \forall x_0 \in \mathcal{X}, \qquad (1)$$

where $x_{t+1} = f(x_t, \pi(x_t))$. This function is the objective we aim to maximize in RL. Note that since we are considering a finite MDP, a state-value function $V^\pi$ can also be viewed as a vector in $\mathbb{R}^{|\mathcal{X}|}$. This perspective can be convenient in the following analysis, so we will interchangeably refer to $V^\pi$ as a function and a vector in the rest of the paper.

Safety is specified through state constraints in the form of an inequality $h(x) \leq 0$, where $h : \mathcal{X} \to \mathbb{R}$ is the constraint function. The constraint function gives an instantaneous safety requirement that should be satisfied at every time step, i.e., we require that

$$h(x_t) \leq 0, \forall t \geq 0. \qquad (2)$$

However, satisfying these infinite-horizon constraints is not always possible. Starting from some states, we may not be able to find a long-term safe policy, and the state constraint will inevitably be violated at some time step in the future. To formally describe and study the long-term safety under state constraints, we introduce the concept of feasibility.

### B. Feasibility in safe reinforcement learning

Feasibility in safe RL includes state feasibility and policy feasibility. State feasibility refers to the satisfiability of the infinite-horizon state constraints, while policy feasibility describes the long-term safety of a specific policy.

*Definition 3.1 (State feasibility):*

1) A state $x_0 \in \mathcal{X}$ is feasible if there exists a policy $\pi$, such that for all $t \geq 0$, we have $h(x_t) \leq 0$, where $x_{t+1} = f(x_t, \pi(x_t))$.
2) A set $\mathrm{X} \subseteq \mathcal{X}$ is a feasible region if all states in this set are feasible.
3) The maximum feasible region $\mathrm{X}^* \subseteq \mathcal{X}$ is the set of all states that are feasible.

By definition of state feasibility, the infinite-horizon constraints (2) can be satisfied if and only if the initial state $x_0$ is in the maximum feasible region. Given a policy in the MDP, some initial states in the maximum feasible region may satisfy the infinite-horizon constraints while others may not. Since safety is our priority concern, we require that the policy we find should render all states in the maximum feasible region safe in the long term. To formally describe this requirement, we introduce the concept of policy feasibility.

*Definition 3.2 (Policy feasibility):*

1) A policy $\pi$ is feasible in a state $x_0 \in \mathcal{X}$ if for any $t \geq 0$, we have $h(x_t) \leq 0$, where $x_{t+1} = f(x_t, \pi(x_t))$.
2) The feasible region of $\pi$, $\mathrm{X}^\pi \subseteq \mathcal{X}$, is the set of all states in which $\pi$ is feasible.

The feasible region of policy is by definition a subset of the maximum feasible region, i.e., for any $\pi$, we have $\mathrm{X}^\pi \subseteq \mathrm{X}^*$. The requirement that a policy renders all states in the maximum feasible region safe in the long term is equivalent

to the fact that the feasible region of this policy equals the maximum feasible region. Another way to understand this requirement is through the working area of a policy. The feasible region of a policy is its safe working area, i.e., the set where it can be safely applied. Due to the safety priority requirement, we want this area to be as large as possible. Furthermore, a larger feasible region may also lead to better control performance because it allows the policy to explore a larger set in the state space, increasing the possibility of obtaining higher rewards.

### C. Objective of safe reinforcement learning

Through the above analysis, we know that the safety requirement in safe RL is to find a policy that has the maximum feasible region. Together with the performance objective of maximizing the state-value function, we arrive at the objective of safe RL:

$$\max_{\pi \in \Pi^*} V^\pi(x), \forall x \in X^*, \tag{3}$$

where $\Pi^*$ is the set of policies with the maximum feasible region, i.e., $\Pi^* = \{\pi | X^\pi = X^*\}$.

## IV. FEASIBILITY FUNCTION AND OPTIMALITY CONDITION

### A. Notation declaration

For simplicity of analysis, we declare the following notation usages.

- Given a non-empty subset $X \subseteq \mathcal{X}$, we define an equivalence relation $=_X$ and a partial order $\leq_X$ on $\mathbb{R}^{|\mathcal{X}|}$:

$$V_1 =_X V_2 \iff V_1(x) = V_2(x), \forall x \in X,$$
$$V_1 \leq_X V_2 \iff V_1(x) \leq V_2(x), \forall x \in X,$$

  where $=$ and $\leq$ are correspondingly the normal equivalence relation and partial order on $\mathbb{R}$. It holds directly from the properties of $=$ and $\leq$ that the definitions of $=_X$ and $\leq_X$ are valid. We will omit the subscript when $X = \mathcal{X}$ since $=_X$ and $\leq_X$ degenerate correspondingly to the normal pointwise equivalence and partial order on $\mathbb{R}^{|\mathcal{X}|}$ under such cases.

- Given a non-empty subset $X \subseteq \mathcal{X}$, we construct a normed space $\left(\mathbb{R}^{|\mathcal{X}|}, =_X, \leq_X, \|\cdot\|_{X,\infty}\right)$ by equipping $\left(\mathbb{R}^{|\mathcal{X}|}, =_X, \leq_X\right)$ with a norm $\|\cdot\|_{X,\infty}$:

$$\|V\|_{X,\infty} = \sup_{x \in X} |V(x)|, \forall V \in \mathbb{R}^{|\mathcal{X}|}.$$

  This is equivalent to the normal infinity norm $\|\cdot\|_\infty$ on $\mathbb{R}^{|\mathcal{X}|}$ when $X = \mathcal{X}$.

- We use $(\cdot)_+$ to denote clipping a function to non-negative values, i.e., $(\cdot)_+ = \max\{\cdot, 0\}$.
- We use $\bar{\cdot}$ to denote the complement of a set, e.g., $\bar{X} = \mathcal{X} \setminus X$.

### B. Feasibility function

Solving problem (3) is challenging because the feasible region of a policy is difficult to represent, and it is even more difficult to decide whether a feasible region equals the maximum one. This difficulty stems from the definition of

feasibility, which involves an infinite number of constraints and is thus almost impossible to check directly. To solve this problem, we introduce a tool called the feasibility function that directly represents a feasible region with its zero-sublevel set. This function also satisfies several desirable properties so that it can be easily computed.

*Definition 4.1 (Feasibility function):* $F^\pi : \mathcal{X} \to \mathbb{R}$ (or $F^\pi \in \mathbb{R}^{|\mathcal{X}|}$ from the vector perspective) is a feasibility function of policy $\pi$ if (1) for any state $x \in \mathcal{X}$, $F^\pi(x) \leq 0 \iff x \in X^\pi$, and (2) there exists an operator $R^\pi : \mathbb{R}^{|\mathcal{X}|} \to \mathbb{R}^{|\mathcal{X}|}$ that satisfies the following three properties:

(a) $R^\pi F^\pi = F^\pi$.
(b) $R^\pi$ is monotonic, i.e., for any $F, F' \in \mathbb{R}^{|\mathcal{X}|}$ such that $F \leq F'$, we have $R^\pi F \leq R^\pi F'$.
(c) $R^\pi$ is a $\gamma$-contraction under the infinity norm, i.e., for any $F, F' \in \mathbb{R}^{|\mathcal{X}|}$, we have $\|R^\pi F - R^\pi F'\|_\infty \leq \gamma \|F - F'\|_\infty$, where $0 < \gamma < 1$ is the discount factor.

The operator $R^\pi$ defined above is called the *risky self-consistency operator*. The contraction property allows us to solve $F^\pi$ by iteratively applying $R^\pi$ from an arbitrary initial function: by Banach's fixed point theorem, for any $F \in \mathbb{R}^{|\mathcal{X}|}$, we have

$$F^\pi = \lim_{k \to \infty} (R^\pi)^k F.$$

To better understand the above definition, we give two examples of feasibility functions.

*Example 4.1 (Cost value function):* The cost value function (CVF) of policy $\pi$, $F^\pi_{\text{CVF}}$, is defined as

$$F^\pi_{\text{CVF}}(x_0) = \sum_{t=0}^\infty \gamma^t c(x_t), \forall x_0 \in \mathcal{X},$$

where $c(x_t) = \mathbf{1}[h(x_t) > 0]$ is the indicator function of constraint violation. The zero-sublevel set of $F^\pi_{\text{CVF}}$ is the feasible region of $\pi$ because

$$F^\pi_{\text{CVF}}(x_0) \leq 0 \iff c(x_t) = 0, \forall t \geq 0$$
$$\iff h(x_t) \leq 0, \forall t \geq 0 \iff x_0 \in X^\pi.$$

The risky self-consistency operator of the CVF $R^\pi_{\text{CVF}}$ is

$$R^\pi_{\text{CVF}} F^\pi_{\text{CVF}}(x) = c(x) + \gamma F^\pi_{\text{CVF}}(f(x, \pi(x))), \forall x \in \mathcal{X}.$$

The CVF is exactly in the same form as the state-value function. Therefore, it inherits all three properties in Definition 4.1 from the state-value function and is thus a feasibility function.

*Example 4.2 (Constraint decay function):* The constraint decay function (CDF) of policy $\pi$, $F^\pi_{\text{CDF}}$, is defined as

$$F^\pi_{\text{CDF}}(x_0) = \gamma^{N^\pi(x_0)}, \forall x_0 \in \mathcal{X},$$

where $N^\pi(x_0)$ is the number of time steps that have elapsed until the first occurrence (if any) of a constraint violation starting from $x_0$ and following $\pi$. If $x_0 \in X^\pi$, $N^\pi(x_0) = +\infty$, and otherwise, $N^\pi(x_0) < \infty$. The zero-sublevel set of $F^\pi_{\text{CDF}}$ is the feasible region of $\pi$ because

$$F^\pi_{\text{CDF}}(x_0) \leq 0 \iff N^\pi(x_0) = +\infty$$
$$\iff h(x_t) \leq 0, \forall t \geq 0 \iff x_0 \in X^\pi.$$

The risky self-consistency operator of the CDF $R^\pi_{\mathrm{CDF}}$ is

$$R^\pi_{\mathrm{CDF}} F^\pi_{\mathrm{CDF}}(x) = c(x) + (1-c(x))\gamma F^\pi_{\mathrm{CDF}}(f(x, \pi(x))), \forall x \in \mathcal{X}.$$

Yang et al. [26] have proved that the CDF satisfies the three properties in Definition 4.1 and is thus a feasibility function.

As shown by the above two examples, the form of a feasibility function is not unique. Any function satisfying the properties of a feasibility function can be used to solve safe RL problems. As we will see in later sections, no matter which form of feasibility function is used, the feasible region will be identical, and the policy will also be the same within the feasible region. This is because we only care about where the feasible states are located, not the values of the feasibility function on these states. Therefore, without loss of generality, we make an assumption that will largely simplify the theoretical analysis of our algorithm: We assume that $F^\pi$ only takes non-negative values, i.e., $F^\pi \geq 0$. Another way to state this assumption is that for any $F^\pi$, we can construct another non-negative feasibility function $F^\pi_+ = \max\{F^\pi, 0\}$. Such a constructed feasibility function will not affect the policy and the state-value function within the feasible region, and thus all theoretical analysis under this assumption also applies to general cases.

As one may notice, the risky self-consistency operator of a feasibility function $R^\pi$ is a counterpart of the *self-consistency operator* of a state-value function $T^\pi : \mathbb{R}^{|\mathcal{X}|} \to \mathbb{R}^{|\mathcal{X}|}$, which is defined as

$$T^\pi V(x) = r(x, \pi(x)) + \gamma V(f(x, \pi(x))), \forall x \in \mathcal{X}. \quad (4)$$

It has been proved that $T^\pi$ also satisfies the three properties in Definition 4.1. Therefore, we can also solve $V^\pi$ by fixed point iteration:

$$V^\pi = \lim_{k \to \infty} (T^\pi)^k V, \forall V \in \mathbb{R}^{|\mathcal{X}|}.$$

## C. Optimality condition for safe reinforcement learning

The ultimate goal of safe RL is to find an optimal policy that maximizes the long-term return while having the maximum feasible region as its feasible region. As a result, all the analysis regarding of optimality in standard RL needs to be extended to involve both value and feasibility. Following the route of standard RL, we start by defining optimal feasibility function, optimal state-value function and optimal policy, and then establish optimality condition which describes the essential properties of them and offers a way to solve them.

Comparison between two feasibility function is direct. One feasibility function $F_1$ is considered to be better than or equal to another feasibility function $F_2$, if $F_1 \leq F_2$. This leads to the optimal feasibility function, which is the infimum of feasibility functions over policies.

*Definition 4.2 (Optimal feasibility function):* $F^* : \mathcal{X} \to \mathbb{R}_{\geq 0}$ is an optimal feasibility function if

$$F^* = \inf_\pi F^\pi. \quad (5)$$

The following proposition links $F^*$ and $X^*$, stating that any policy having $F^*$ as its feasibility function must also have $X^*$ as its feasible region.

*Proposition 4.1:* For any policy $\pi$ such that $F^\pi = F^*$, we have $X^\pi = X^*$.

With feasibility under consideration, a state-value function $V^\pi$ is no longer meaningful at every state, but only within the corresponding feasible region $X^\pi$. Consequently, comparison between two state-value function is generally restricted to a certain feasible region. Besides, since feasibility should be a prior consideration to return, the optimal state-value function is defined as the supremum of state-value functions over policies with the maximum feasible region, instead of all policies.

*Definition 4.3 (Optimal state-value function):* $V^* : \mathcal{X} \to \mathbb{R}$ is an optimal state-value function if

$$V^* =_{X^*} \sup_{\pi \in \Pi^*} V^\pi, \quad (6)$$

where $\Pi^* = \{\pi | X^\pi = X^*\}$. Note that the supremum here is based on partial order $\leq_{X^*}$, as can be inferred from the use of $=_{X^*}$ in definition.

An optimal policy in safe RL is defined on top of $F^*$ and $V^*$ as follows.

*Definition 4.4 (Optimal policy):* $\pi^* : \mathcal{X} \to \mathcal{U}$ is an optimal policy if

$$F^{\pi^*} = F^*, \quad (7a)$$

$$V^{\pi^*} =_{X^*} V^*. \quad (7b)$$

In order to solve the above-defined optimal functions and policy, we need also optimality version of self-consistency operators, which are termed *risky Bellman operator* and *feasible Bellman operator*, for feasibility functions and state-value functions correspondingly.

*Definition 4.5 (Risky Bellman operator):* The risky Bellman operator $R^* : \mathbb{R}^{|\mathcal{X}|} \to \mathbb{R}^{|\mathcal{X}|}$ is defined as

$$R^* F = \inf_\pi R^\pi F. \quad (8)$$

*Definition 4.6 (Feasible Bellman operator):* The feasible Bellman operator $T^* : \mathbb{R}^{|\mathcal{X}|} \to \mathbb{R}^{|\mathcal{X}|}$ is defined as

$$T^* V(x) = \begin{cases} \sup_{\pi \in \Pi^*} T^\pi V(x) & x \in X^* \\ V(x) & x \notin X^* \end{cases} \quad (9)$$

Until now, we have defined everything we want to solve (i.e., $F^*$, $V^*$ and $\pi^*$) and tools for solving them (i.e., $R^*$ and $T^*$). The following important theorem theoretically justifies the iterative use of risky Bellman operator and feasible Bellman operator for solving $F^*$ and $V^*$, correspondingly. It also shows the way of obtaining an optimal policy from the solved optimal functions.

*Theorem 4.1 (Optimality condition):* Let $F^*$ be an optimal feasibility function and $V^*$ be an optimal state-value function, we have

1) Any policy $\pi$ that satisfies

$$R^\pi F^* = R^* F^*, \quad (10a)$$

$$T^\pi V^* =_{X^*} T^* V^*, \quad (10b)$$

   is an optimal policy.

2) It holds that

$$F^* = R^* F^*, \quad (11a)$$

$$V^* = T^* V^*. \quad (11b)$$

*Proof:* For any policy $\pi$, we have

$$F^\pi \geq F^*$$
$$\implies \quad F^\pi = R^\pi F^\pi \geq R^\pi F^* \qquad\qquad \text{(Def. 4.1)}$$
$$\implies \quad \inf_\pi F^\pi \geq \inf_\pi R^\pi F^*$$
$$\implies \quad F^* \geq R^* F^* \qquad\qquad \text{(Def. 4.2\&4.5)}$$
$$\implies \quad F^* \geq R^\pi F^* \qquad\qquad \text{(By Eq. (10a))}$$
$$\implies \quad F^* \geq (R^\pi)^k F^* \qquad\qquad \text{(Def. 4.1)}$$
$$\implies \quad F^* \geq F^\pi \qquad\qquad \text{(Let } k \to \infty)$$
$$\implies \quad F^* = F^\pi, \qquad\qquad \text{(Def. 4.2)}$$

which proves the first condition (7a) of an optimal policy. Furthermore, we have

$$F^* = F^\pi = R^\pi F^\pi = R^\pi F^* = R^* F^*,$$

which proves Eq. (11a).

Since $F^* = F^\pi$, we have $X^\pi = X^*$ by Proposition 4.1, which means $\pi \in \Pi^*$. Thus, we have

$$V^\pi \leq \sup_{\pi \in \Pi^*} V^\pi*$$
$$\implies \quad V^\pi \leq_{X^*} V^* \qquad\qquad \text{(Def. 4.3)}$$
$$\implies \quad V^\pi = T^\pi V^\pi \leq_{X^*} T^\pi V^* \qquad\qquad \text{(Monotonicity)}$$
$$\implies \quad \sup_{\pi \in \Pi^*} V^\pi \leq_{X^*} \sup_{\pi \in \Pi^*} T^\pi V^*$$
$$\implies \quad V^* \leq_{X^*} T^* V^* \qquad\qquad \text{(Def. 4.3\&4.6)}$$
$$\implies \quad V^* \leq_{X^*} T^\pi V^* \qquad\qquad \text{(By Eq. (10b))}$$
$$\implies \quad V^* \leq_{X^*} (T^\pi)^k V^* \qquad\qquad \text{(Monotonicity)}$$
$$\implies \quad V^* \leq_{X^*} V^\pi \qquad\qquad \text{(Let } k \to \infty)$$
$$\implies \quad V^* =_{X^*} V^\pi, \qquad\qquad \text{(Def. 4.3)}$$

which proves the second condition (7b) of an optimal policy. Furthermore, we have

$$V^* =_{X^*} V^\pi = T^\pi V^\pi =_{X^*} T^\pi V^* =_{X^*} T^* V^*,$$

which, together with $V^* =_{\bar{X}^*} T^* V^*$, proves Eq. (11b) and thus finishes the proof of the theorem. ■

Eq. (11a) is called the *risky Bellman equation*, and Eq. (11b) is called the *feasible Bellman equation*. They indicate that, by iteratively applying $R^*$ and $T^*$, we will eventually converge to fixed points of them, which are $F^*$ and $V^*$, correspondingly. And item 1 of Theorem 4.1 further states that, if we behave "greedily" (in terms of both feasibility and return) according to the solved $F^*$ and $V^*$, we can obtain an optimal policy $\pi^*$.

As suggested by their names, all the definitions and theorem introduced in this chapter together form a safe-RL extension of the optimality theory in standard RL, which is a long missing yet of significant importance piece of safe RL.

## V. FEASIBLE POLICY ITERATION

Through the above analysis, we know that the essence of safe RL is to solve the two Bellman equations (11a) and (11b). To achieve this goal, we propose an iterative algorithm called FPI, which alternates between three steps: (1) PEV, (2) RID, and (3) region-wise PIM. PEV is the same as in standard RL, which computes the state-value function of the current policy. RID computes the feasibility function, which gives the feasible region of the current policy. Region-wise PIM updates the policy to maximize its state-value function and minimize its feasibility function simultaneously.

### A. Policy evaluation

In PEV, we solve the self-consistency condition of the state-value function:

$$V^\pi(x) = r(x, \pi(x)) + \gamma V^\pi(f(x, \pi(x))), \forall x \in \mathcal{X}. \quad (12)$$

As previously mentioned, this equation can be solved by iteratively applying the self-consistency operator $T^\pi$ to an arbitrary initial state-value function. By the contraction property of $T^\pi$, we can expect geometric convergence to $V^\pi$. However, in finite state space, we have a more straightforward method to solve $V^\pi$ because Eq. (12) has an analytical solution. To see this, first construct a state transition matrix $P^\pi \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$, where $P_{ij}^\pi$ represents the probability of transferring from $x_i$ to $x_j$ under policy $\pi$. In a deterministic MDP, each row of $P$ has a single element that equals one, which corresponds to the unique next state, and all other elements are zero. Therefore, we have $P_{ij}^\pi = \mathbf{1}[f(x_i, \pi(x_i)) = x_j]$. Using this state transition matrix, and with a slight abuse of notation, we can write Eq. (12) as

$$V^\pi = r^\pi + \gamma P^\pi V^\pi, \quad (13)$$

where here we view $V^\pi$ as a vector with $|\mathcal{X}|$ elements and $r^\pi$ also a vector with elements $r^\pi(x) = r(x, \pi(x))$. With a rearrangement of Eq. (13) and noticing that the matrix $I - \gamma P^\pi$ is invertible, we can readily obtain the solution as $V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$.

### B. Region identification

In RID, we solve the risky self-consistency condition of the feasibility function:

$$F^\pi(x) = R^\pi F^\pi(x), \forall x \in \mathcal{X}. \quad (14)$$

Similar to the self-consistency condition, Eq. (14) can also be solved by iteratively applying the risky self-consistency operator $R^\pi$, which provides a geometric convergence. Comparing Eq. (14) with Eq. (12), a question is can we also analytically solve Eq. (14)? Without the explicit form of $R^\pi$, we cannot guarantee an analytical solution exists. Fortunately, an analytical solution does exist for many widely-used feasibility functions. The CVF in Example 4.1 obviously can be analytically solved because it has the same form as the state-value function. As a non-trivial example, we show that the CDF in Example 4.2 can also be analytically solved. The risky self-consistency condition of the CDF is

$$F^\pi(x) = c(x) + (1 - c(x))\gamma F^\pi(f(x, \pi(x))), \forall x \in \mathcal{X},$$

which can be written as

$$F^\pi = c + \gamma C P^\pi F^\pi,$$

where $C \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{X}|}$ is a diagonal matrix with its elements $C_{ii} = 1 - c_i$. The analytical solution of this equation can be readily obtained as $F^\pi = (I - \gamma C P^\pi)^{-1} c$.

## C. Region-wise policy improvement

In region-wise PIM, we separately update the policy on states inside and outside the current feasible region. Let $\pi$ be the original policy and $\pi'$ be the new policy. For states outside the feasible region, we update the policy to minimize the feasibility function:

$$\pi'(x) = \arg\inf_{\tilde{\pi}} R^{\tilde{\pi}} F^\pi(x), \forall x \in \bar{\mathrm{X}}^\pi. \tag{15}$$

Using the risky Bellman operator, this update rule can be equivalently written as

$$R^{\pi'} F^\pi =_{\bar{\mathrm{X}}^\pi} R^* F^\pi.$$

Since we assume that $F^\pi$ equals zero inside the feasible region, the above equation can be simplified to $R^{\pi'} F^\pi = R^* F^\pi$. For states inside the feasible region, we update the policy to maximize the state-value function, but the new policy should keep the current feasible region forward invariant:

$$\pi'(x) = \arg\sup_{\tilde{\pi} \in \Pi^\pi} T^{\tilde{\pi}} V^\pi(x), \forall x \in \mathrm{X}^\pi, \tag{16}$$

where $\Pi^\pi$ is the set of policies that keep $\mathrm{X}^\pi$ forward invariant, i.e.,

$$\Pi^\pi = \{\tilde{\pi} | \forall x \in \mathrm{X}^\pi, f(x, \tilde{\pi}(x)) \in \mathrm{X}^\pi\}.$$

To facilitate theoretical analysis of our algorithm, we introduce an operator to describe this update rule.

*Definition 5.1 (Restricted feasible Bellman operator):* The restricted feasible Bellman operator on the feasible region of $\pi$, $T_\pi^* : \mathbb{R}^{|\mathcal{X}|} \to \mathbb{R}^{|\mathcal{X}|}$ is defined as

$$T_\pi^* V(x) = \begin{cases} \sup_{\tilde{\pi} \in \Pi^\pi} T^{\tilde{\pi}} V(x) & x \in \mathrm{X}^\pi, \\ V(x) & x \notin \mathrm{X}^\pi. \end{cases} \tag{17}$$

Using the restricted feasible Bellman operator, the update rule (16) can be equivalently written as

$$T^{\pi'} V^\pi =_{\mathrm{X}^\pi} T_\pi^* V^\pi.$$

## D. Theoretical analysis

The theoretical significance of FPI lies in its two fundamental guarantees: monotonic improvement and geometric convergence, which are analogous to those in PI, but are largely extended in the context of safe RL. Monotonic improvement refers to the consistent decrease of the feasibility function and the consistent increase of the state-value function. Similarly, geometric convergence applies to both the feasibility and state-value functions. Despite the theoretical foundation established by PI, extending these properties to FPI is far from straightforward. A key challenge is that the Bellman operator is restricted to a feasible region during PIM. This restriction ties the monotonicity of the state-value function to that of the feasible region and complicates the convergence mechanism of the state-value function. To begin the proof, we first define the restricted optimal state-value function on a feasible region.

*Definition 5.2 (Restricted optimal state-value function):* $V_\pi^* : \mathcal{X} \to \mathbb{R}$ is a restricted optimal state-value function on the feasible region of $\pi$ if

$$V_\pi^* =_{\mathrm{X}^\pi} \sup_{\pi' \in \Pi^\pi} V^{\pi'}. \tag{18}$$

The restricted optimal state-value function is different from, or more precisely, lower than the optimal state-value function because it maximizes over a subset of the policy space rather than the entire space. This function represents the best achievable solution when performing PIM inside a given feasible region. A notable special case arises when the feasible region equals the maximum one. In this case, the restricted optimal state-value function becomes the optimal state-value function.

*Proposition 5.1:* For any policy $\pi$ such that $\mathrm{X}^\pi = \mathrm{X}^*$, $V_\pi^*$ is an optimal state-value function.

*Proof:* If we can prove $\Pi^\pi = \Pi^*$, then we have

$$V_\pi^* =_{\mathrm{X}^*} \sup_{\pi' \in \Pi^*} V^{\pi'},$$

which means $V_\pi^*$ is an optimal state-value function. To prove $\Pi^\pi = \Pi^*$, we show that $\Pi^\pi \subseteq \Pi^*$ and $\Pi^* \subseteq \Pi^\pi$. First, for any policy $\pi' \in \Pi^\pi$, we have

$$\forall x \in \mathrm{X}^\pi = \mathrm{X}^*, f(x, \pi'(x)) \in \mathrm{X}^\pi = \mathrm{X}^*.$$

Thus, $\mathrm{X}^*$ is forward invariant under $\pi'$, which mean $\mathrm{X}^* \subseteq \mathrm{X}^{\pi'}$. Together with $\mathrm{X}^{\pi'} \subseteq \mathrm{X}^*$, we have $\mathrm{X}^{\pi'} = \mathrm{X}^*$, and thus $\pi' \in \Pi^*$. By the arbitrariness of $\pi'$, we conclude that $\Pi^\pi \subseteq \Pi^*$. On the other hand, for any policy $\pi' \in \Pi^*$, using that $\mathrm{X}^{\pi'}$ is forward invariant under $\pi'$, we have

$$\forall x \in \mathrm{X}^{\pi'} = \mathrm{X}^* = \mathrm{X}^\pi, f(x, \pi'(x)) \in \mathrm{X}^{\pi'} = \mathrm{X}^* = \mathrm{X}^\pi.$$

Thus, we have $\pi' \in \Pi^\pi$, which proves that $\Pi^* \subseteq \Pi^\pi$ and, therefore, $\Pi^\pi = \Pi^*$. ∎

This proposition reveals the dependencies of reward performance on the size of the feasible region: achieving optimal reward performance requires obtaining the maximum feasible region. This is because the maximum feasible region provides the largest policy space for searching the optimal policy. Just as the optimality condition for the unrestricted policy space, the restricted optimal state-value function follows a similar condition related to the restricted feasible Bellman operator.

*Lemma 5.1 (Restricted optimality condition):* Let $V_\pi^*$ be a restricted optimal state-value function on the feasible region of $\pi$, we have

$$V_\pi^* = T_\pi^* V_\pi^*. \tag{19}$$

*Proof:* For any policy $\pi' \in \Pi^\pi$, we have

$$\begin{aligned} & V^{\pi'} \le \sup_{\tilde{\pi} \in \Pi^\pi} V^{\tilde{\pi}} =_{\mathrm{X}^\pi} V_\pi^* \\ \implies & T^{\pi'} V^{\pi'} \le_{\mathrm{X}^\pi} T^{\pi'} V_\pi^* \qquad (\mathrm{X}^\pi \subseteq \mathrm{X}^{\pi'}) \\ \implies & V^{\pi'} \le_{\mathrm{X}^\pi} T^{\pi'} V_\pi^* \\ \implies & \sup_{\pi' \in \Pi^\pi} V^{\pi'} \le_{\mathrm{X}^\pi} \sup_{\pi' \in \Pi^\pi} T^{\pi'} V_\pi^* \\ \implies & V_\pi^* \le_{\mathrm{X}^\pi} T_\pi^* V_\pi^*. \end{aligned} \tag{20}$$

Let $\pi'^* = \arg\sup_{\pi' \in \Pi^\pi} T^{\pi'} V_\pi^*$, then we have

$$T^{\pi'^*} V_\pi^* =_{\mathrm{X}^\pi} T_\pi^* V_\pi^*.$$

Chaining this with Eq. (20), we have

$$\begin{aligned} & V_\pi^* \le_{\mathrm{X}^\pi} T^{\pi'^*} V_\pi^* \\ \implies & V_\pi^* \le_{\mathrm{X}^\pi} (T^{\pi'^*})^k V_\pi^* \qquad \text{(Monotonicity)} \\ \implies & V_\pi^* \le_{\mathrm{X}^\pi} V^{\pi'^*}. \qquad \text{(Let } k \to \infty) \end{aligned}$$

Together with $V_\pi^* \geq_{X^\pi} V^{\pi'^*}$, we have $V_\pi^* =_{X^\pi} V^{\pi'^*}$. Therefore, we have

$$V_\pi^* =_{X^\pi} V^{\pi'^*} = T^{\pi'^*} V^{\pi'^*} =_{X^\pi} T^{\pi'^*} V_\pi^* =_{X^\pi} T_\pi^* V_\pi^*.$$

Since $T_\pi^* V_\pi^* =_{\bar{X}^\pi} V_\pi^*$, we conclude that $V_\pi^* = T_\pi^* V_\pi^*$. ■

Eq. (19) is called the *restricted feasible Bellman equation*. Through the iterations of FPI, we are essentially solving this equation to obtain the restricted optimal state-value function. As we will later illustrate, this equation also enables us to solve the value function through an iterative process with convergence guarantees. A key factor enabling this convergence is the monotonicity of the Bellman operators.

*Lemma 5.2 (Monotonicity of operators):* For any policy $\pi$, set $X \subseteq X^\pi$, and functions $F, F' \in \mathbb{R}^{|\mathcal{X}|}, V, V' \in \mathbb{R}^{|\mathcal{X}|}$ such that $F \leq F', V \leq_X V'$, we have

$$R^* F \leq R^* F', \tag{21a}$$

$$T_\pi^* V \leq_X T_\pi^* V'. \tag{21b}$$

*Proof:* By monotonicity of $R^{\pi'}$, we have

$$R^{\pi'} F \leq R^{\pi'} F'$$
$$\implies \quad \inf_{\pi'} R^{\pi'} F \leq \inf_{\pi'} R^{\pi'} F'$$
$$\implies \quad R^* F \leq R^* F', \tag{Def. 4.5}$$

which proves Eq. (21a). By monotonicity of $T^{\pi'}$, we have

$$T^{\pi'} V \leq_X T^{\pi'} V'$$
$$\implies \quad \sup_{\pi' \in \Pi^\pi} T^{\pi'} V \leq_X \sup_{\pi' \in \Pi^\pi} T^{\pi'} V'$$
$$\implies \quad T_\pi^* V \leq_X T_\pi^* V', \tag{Def. 4.5}$$

which proves Eq. (21b). ■

Another important property that enables the iterative solving of the Bellman equations is the contraction property. This property is the core reason why FPI can achieve a geometric convergence speed.

*Lemma 5.3 ($\gamma$-contraction of Bellman operators):* For any policy $\pi$, any set $X \subseteq X^\pi$, and any functions $F, F' \in \mathbb{R}^{|\mathcal{X}|}, V, V' \in \mathbb{R}^{|\mathcal{X}|}$, we have

$$\|R^* F - R^* F'\|_\infty \leq \gamma \|F - F'\|_\infty, \tag{22a}$$

$$\|T_\pi^* V - T_\pi^* V'\|_{X,\infty} \leq \gamma \|V - V'\|_{X,\infty}. \tag{22b}$$

*Proof:* We first prove that

$$\|\inf_\pi R^\pi F - \inf_\pi R^\pi F'\|_\infty \leq \sup_\pi \|R^\pi F - R^\pi F'\|_\infty.$$

Let $\pi_1 = \arg\inf_\pi R^\pi F, \pi_2 = \arg\inf_\pi R^\pi F'$, then we have

$$\|\inf_\pi R^\pi F - \inf_\pi R^\pi F'\|_\infty = \|R^{\pi_1} F - R^{\pi_2} F'\|_\infty$$
$$\leq \max\{\|R^{\pi_1} F - R^{\pi_1} F'\|_\infty, \|R^{\pi_2} F - R^{\pi_2} F'\|_\infty\}$$
$$\leq \sup_\pi \|R^\pi F - R^\pi F'\|_\infty.$$

Therefore,

$$\|R^* F - R^* F'\|_\infty = \|\inf_\pi R^\pi F - \inf_\pi R^\pi F'\|_\infty$$
$$\leq \sup_\pi \|R^\pi F - R^\pi F'\|_\infty$$
$$\leq \gamma \|F - F'\|_\infty,$$

which proves Eq. (22a). In the last inequality, we used the $\gamma$-contraction property of $R^\pi$.

With similar reasoning as above, we have

$$\|\sup_{\pi' \in \Pi^\pi} T^{\pi'} V - \sup_{\pi' \in \Pi^\pi} T^{\pi'} V'\|_{X,\infty} \leq \sup_{\pi' \in \Pi^\pi} \|T^{\pi'} V - T^{\pi'} V'\|_{X,\infty}.$$

By $T_\pi^* V =_{X^\pi} \sup_{\pi' \in \Pi^\pi} T^{\pi'} V$ and $X \subseteq X^\pi$, we have

$$T_\pi^* V =_X \sup_{\pi' \in \Pi^\pi} T^{\pi'} V.$$

Therefore,

$$\|T_\pi^* V - T_\pi^* V'\|_{X,\infty} = \|\sup_{\pi' \in \Pi^\pi} T^{\pi'} V - \sup_{\pi' \in \Pi^\pi} T^{\pi'} V'\|_{X,\infty}$$
$$\leq \sup_{\pi' \in \Pi^\pi} \|T^{\pi'} V - T^{\pi'} V'\|_{X,\infty}$$
$$\leq \gamma \|V - V'\|_{X,\infty},$$

which proves Eq. (22b). In the last inequality, we used the $\gamma$-contraction property of $T^{\pi'}$. ■

With the above lemmas regarding the Bellman operators, we are ready to introduce the first key theoretical property of FPI: monotonic improvement. This property characterizes the relationship between two policies generated in successive iterations of FPI.

*Theorem 5.1 (Monotonic improvement):* Let $\pi, \pi'$ be two policies such that

$$R^{\pi'} F^\pi = R^* F^\pi, \tag{23a}$$

$$T^{\pi'} V^\pi =_{X^\pi} T_\pi^* V^\pi. \tag{23b}$$

Then, we have

$$F^\pi \geq R^* F^\pi \geq F^{\pi'}, \tag{24a}$$

$$V^\pi \leq T_\pi^* V^\pi \leq_{X^\pi} V^{\pi'}. \tag{24b}$$

*Proof:* By $F^\pi = R^\pi F^\pi \geq R^* F^\pi$, we have

$$F^\pi \geq R^* F^\pi = R^{\pi'} F^\pi. \tag{25}$$

We prove by induction that for $k \geq 1$,

$$F^\pi \geq R^* F^\pi \geq (R^{\pi'})^k F^\pi. \tag{26}$$

From this, Eq. (24a) will follow by taking $k \to \infty$ on both sides. The case of $k = 1$ has already been proved. Assume that the required inequality holds for $k \geq 1$, and we prove that it also holds for $k+1$. Applying $R^{\pi'}$ on both sides of Eq. (26) and using monotonicity of $R^{\pi'}$, we have

$$R^{\pi'} F^\pi \geq (R^{\pi'})^{k+1} F^\pi.$$

Chaining this with Eq. (25), we have

$$F^\pi \geq R^* F^\pi = R^{\pi'} F^\pi \geq (R^{\pi'})^{k+1} F^\pi,$$

which finishes the inductive step, and thus the proof of Eq. (26).

By $V^\pi = T^\pi V^\pi \leq T_\pi^* V^\pi$, we have

$$V^\pi \leq T_\pi^* V^\pi =_{X^\pi} T^{\pi'} V^\pi. \tag{27}$$

We prove by induction that for $k \geq 1$,

$$V^\pi \leq T_\pi^* V^\pi \leq_{X^\pi} (T^{\pi'})^k V^\pi. \tag{28}$$

From this, Eq. (24b) will follow by taking $k \to \infty$ on both sides. The case of $k = 1$ has already been proved. Assume that the required inequality holds for $k \geq 1$, and we prove that

it also holds for $k+1$. Applying $T^{\pi'}$ on both sides of Eq. (28) and using monotonicity of $T^{\pi'}$, we have

$$T^{\pi'}V^{\pi} \leq_{X^{\pi}} (T^{\pi'})^{k+1}V^{\pi}.$$

Chaining this with Eq. (27), we have

$$V^{\pi} \leq T^*_{\pi}V^{\pi} =_{X^{\pi}} T^{\pi'}V^{\pi} \leq_{X^{\pi}} (T^{\pi'})^{k+1}V^{\pi},$$

which finishes the inductive step, and thus the proof of Eq. (28). ■

Monotonic improvement ensures that the policy generated in successive iterations of FPI surpasses its predecessors in terms of both safety and performance. This property enables FPI to avoid the oscillation problem commonly observed in other methods, such as the Lagrange method, and as we will see in our experiments, this leads to stable training processes. A direct proposition of the monotonic decrease of the feasibility function is that the feasible region monotonically expands.

*Proposition 5.2 (Monotonic region expansion):* Let $\pi, \pi'$ be two policies such that $R^{\pi'}F^{\pi} = R^*F^{\pi}$, we have $X^{\pi} \subseteq X^{\pi'}$.

*Proof:* By Theorem 5.1, we have $F^{\pi} \geq F^{\pi'}$. Thus, for any $x \in X^{\pi}$, we have $F^{\pi'}(x) \leq F^{\pi}(x) \leq 0$, which means $x \in X^{\pi'}$. Therefore, $X^{\pi} \subseteq X^{\pi'}$. ■

Before proving geometric convergence, we need to answer another question: when optimizing the state-value function within different feasible regions, do larger regions always yield higher values? This question is crucial because the feasible regions in each iteration of FPI are different, but the monotonicity of the state-value function we just proved is based on a fixed region. To extend the monotonicity to these changing regions, we need to show that optimizing over a sequence of expanding regions still preserves monotonicity, which is stated by the following lemma.

*Lemma 5.4 (Feasible region superiority):* Let $\pi, \pi'$ be two policies such that $X^{\pi} \subseteq X^{\pi'}$, we have

$$T^*_{\pi}V \leq_{X^{\pi}} T^*_{\pi'}V. \tag{29}$$

*Proof:* We first prove that for any policy $\tilde{\pi} \in \Pi^{\pi}$, there exists a policy $\tilde{\pi}' \in \Pi^{\pi'}$ such that $\tilde{\pi}' =_{X^{\pi}} \tilde{\pi}$. Such a policy can be constructed as

$$\tilde{\pi}'(x) = \begin{cases} \tilde{\pi}(x) & x \in X^{\pi}, \\ \pi'(x) & x \notin X^{\pi}, \end{cases}$$

which naturally satisfies $\tilde{\pi}' =_{X^{\pi}} \tilde{\pi}$, and the left is to prove $\tilde{\pi}' \in \Pi^{\pi'}$. To this end, we prove that for any $x \in X^{\pi'}$, we have $f(x, \tilde{\pi}'(x)) \in X^{\pi'}$ by considering two cases: (1) $x \in X^{\pi}$, and (2) $x \in X^{\pi'} \setminus X^{\pi}$. If $x \in X^{\pi}$, we have

$$f(x, \tilde{\pi}'(x)) = f(x, \tilde{\pi}(x)) \in X^{\pi} \subseteq X^{\pi'}.$$

If $x \in X^{\pi'} \setminus X^{\pi}$, we have

$$f(x, \tilde{\pi}'(x)) = f(x, \pi'(x)) \in X^{\pi'}.$$

Thus, we conclude that $\tilde{\pi}' \in \Pi^{\pi'}$. Let $\tilde{\pi}^* = \arg\sup_{\tilde{\pi} \in \Pi^{\pi}} T^{\tilde{\pi}}V$. By the conclusion above, there exists $\tilde{\pi}'^* \in \Pi^{\pi'}$ such that $\tilde{\pi}'^* =_{X^{\pi}} \tilde{\pi}^*$. Thus, we have

$$T^*_{\pi}V =_{X^{\pi}} \sup_{\tilde{\pi} \in \Pi^{\pi}} T^{\tilde{\pi}}V =_{X^{\pi}} T^{\tilde{\pi}^*}V$$
$$=_{X^{\pi}} T^{\tilde{\pi}'^*}V \leq_{X^{\pi}} \sup_{\tilde{\pi}' \in \Pi^{\pi'}} T^{\tilde{\pi}'}V =_{X^{\pi}} T^*_{\pi'}V,$$

where the last equality uses that $X^{\pi} \subseteq X^{\pi'}$. ■

Now we have everything needed for proving the most important property of FPI: geometric convergence. This property includes convergence of the feasibility function to the global optimum and convergence of the state-value function to the restricted optimum on the feasible regions.

*Theorem 5.2 (Geometric convergence):* Let $\{\pi_k\}_{k \geq 0}$ be the sequence of policies produced by FPI. Then, for any $k \geq 0$,

$$\|F^* - F^{\pi_k}\|_{\infty} \leq \gamma^k \|F^* - F^{\pi_0}\|_{\infty}, \tag{30a}$$

$$\|(V^*_{\pi_0} - V^{\pi_k})_+\|_{X^{\pi_0},\infty} \leq \gamma^k \|V^*_{\pi_0} - V^{\pi_0}\|_{X^{\pi_0},\infty}. \tag{30b}$$

*Proof:* By Theorem 5.1, we have

$$R^*F^{\pi_0} \geq F^{\pi_1}, \tag{31a}$$

$$T^*_{\pi_0}V^{\pi_0} \leq_{X^{\pi_0}} V^{\pi_1}. \tag{31b}$$

Applying $R^*$ to both sides of Eq. (31a) and using the monotonicity of $R^*$, we have

$$(R^*)^2 F^{\pi_0} \geq R^* F^{\pi_1} \geq F^{\pi_2},$$

where the second inequality is obtained by using Theorem 5.1 again. With the same reasoning, we have

$$(R^*)^k F^{\pi_0} \geq F^{\pi_k}.$$

Thus,

$$0 \leq F^{\pi_k} - F^* \leq (R^*)^k F^{\pi_0} - F^*.$$

Taking the infinity norm, we have

$$\|F^{\pi_k} - F^*\|_{\infty} \leq \|(R^*)^k F^{\pi_0} - F^*\|_{\infty}$$
$$= \|(R^*)^k F^{\pi_0} - (R^*)^k F^*\|_{\infty}$$
$$\leq \gamma^k \|F^{\pi_0} - F^*\|_{\infty},$$

which proves Eq. (30a). In the equality above we used the optimality condition and in the last inequality we used the contraction property.

Applying $T^*_{\pi_1}$ to both sides of Eq. (31b) and using the monotonicity of $T^*_{\pi_1}$, we have

$$T^*_{\pi_1}T^*_{\pi_0}V^{\pi_0} \leq_{X^{\pi_0}} T^*_{\pi_1}V^{\pi_1} \leq_{X^{\pi_1}} V^{\pi_2}.$$

By Lemma 5.4 and $X^{\pi_0} \subseteq X^{\pi_1}$, we have

$$(T^*_{\pi_0})^2 V^{\pi_0} \leq_{X^{\pi_0}} T^*_{\pi_1}T^*_{\pi_0}V^{\pi_0} \leq_{X^{\pi_0}} V^{\pi_2}.$$

With the same reasoning, we have

$$(T^*_{\pi_0})^k V^{\pi_0} \leq_{X^{\pi_0}} V^{\pi_k}.$$

Thus,

$$V^*_{\pi_0} - V^{\pi_k} \leq_{X^{\pi_0}} V^*_{\pi_0} - (T^*_{\pi_0})^k V^{\pi_0}.$$

Applying $(\cdot)_+$ on the left hand side and taking the absolute value of the right hand side, we have

$$(V^*_{\pi_0} - V^{\pi_k})_+ \leq_{X^{\pi_0}} |V^*_{\pi_0} - (T^*_{\pi_0})^k V^{\pi_0}|.$$

Taking the infinity norm on both sides, we have

$$\|(V^*_{\pi_0} - V^{\pi_k})_+\|_{X^{\pi_0},\infty} \leq \|V^*_{\pi_0} - (T^*_{\pi_0})^k V^{\pi_0}\|_{X^{\pi_0},\infty}$$
$$= \|(T^*_{\pi_0})^k V^*_{\pi_0} - (T^*_{\pi_0})^k V^{\pi_0}\|_{X^{\pi_0},\infty}$$
$$\leq \gamma^k \|V^*_{\pi_0} - V^{\pi_0}\|_{X^{\pi_0},\infty},$$

where the equality uses the restricted optimality condition. This proves Eq. (30b) and thus finished the proof. ∎

As stated in this theorem, the optimal state-value function is restricted to the feasible region of the initial policy. While this may seem overly restrictive, it is important to note that any policy in the sequence $\{\pi_k\}_{k\geq0}$ can be viewed as an initial policy. This is because any subsequence of policies is still generated by FPI and retains the convergence property. Given the geometric convergence of the feasibility function, we can expect that the feasible region to quickly approach the maximum feasible region. When the feasible region grows to a satisfying size, we can analyze the policy sequence from this point onward, leveraging the convergence of the state-value function. In this case, the state-value function continues to exhibit geometric convergence to the optimum restricted to this feasible region, which is already very close to the global optimum.

In MDPs with finite state and action spaces, FPI achieves convergence within a finite number of iterations. This is because the total number of deterministic policies is finite, and both the feasibility function and state-value function improve monotonically. If two policies share identical feasibility functions, the maximum feasible region has been reached. Subsequently, if the state-value functions of any two successive policies are identical within this region, the optimal state-value function has been attained. However, this analysis is naive because it first requires the feasibility function to converge before turning to the state-value function. This is impractical because of the low computational efficiency. The significance of Theorem 5.2 is to show that state-value computation can start at any point, even at the very beginning, instead of waiting for the feasibility function to converge. Regardless of the starting point, the state-value function always exhibits geometric convergence to the optimal value, restricted to a feasible region that itself converges to the maximum feasible region.

## VI. PRACTICAL IMPLEMENTATIONS

Our theoretical analysis of FPI is based on finite state and action spaces, where we can use tables to store values of feasibility function, state-value function, and policy and update them by exact value assignments. However, when it comes to infinite state and action spaces, tabular methods no longer apply, and function approximation becomes necessary. Similar to standard deep RL algorithms, we approximate the feasibility function, state-value function, and policy using neural networks and update them via gradient descent. In this section, we introduce some practical implementations for function approximation used in our experiments.

On the basis of FPI, we adopt SAC [9] for reward optimization and denote the resulting algorithm as FPI-SAC. Our algorithm learns an action feasibility function $G_\phi(x, u)$ represented by a neural network with parameter $\phi$, which is a Q-function-like variant of the CDF in Example 4.2. It takes the current state and action as input and outputs the CDF value of the next state, i.e., $G^\pi(x, u) = F^\pi(f(x, u))$. FPI-SAC also learns two action-value networks $Q_{\omega_1}, Q_{\omega_2}$ and a stochastic policy network $\pi_\theta$.

We use a sigmoid output activation function for $G_\phi$ because its value is within $[0, 1]$ by definition. Incorporating this prior structure into the network makes it easier to converge because it eliminates excessively large errors introduced by bootstrapping. To avoid the gradient vanishing problem caused by sigmoid function, we use a cross-entropy loss to train $G_\phi$, with the right-hand side of (14) as the training label:

$$L_G(\phi) = -\mathbb{E}_{(x,u)\sim\mathcal{D}}\{y_G(x)\log G_\phi(x, u)$$
$$+(1 - y_G(x))\log(1 - G_\phi(x, u))\}, \quad (32)$$
$$y_G(x) = c(x) + (1 - c(x))\gamma G_{\bar{\phi}}(x', u'),$$

where $\mathcal{D}$ is the replay buffer and $\bar{\phi}$ is the parameters of the target network for $G_\phi$, which is updated at a slower rate to stabilize the training process.

The loss functions of the Q networks are the same as those in SAC:

$$L_Q(\omega_i) = \mathbb{E}_{(x,u,r,x')\sim\mathcal{D}}\{(y_Q - Q_{\omega_i}(x, u))^2\},$$
$$y_Q = r + \gamma(\min_{j\in\{1,2\}} Q_{\bar{\omega}_j}(x', u') - \alpha\log\pi_\theta(u'|x')),$$
$$(33)$$

where $i \in \{1, 2\}$, $\bar{\omega}_j$ are the parameters of the target Q networks, and $\alpha$ is the temperature.

The loss function of the policy is computed separately inside and outside the feasible region. We introduce a feasibility threshold $0 < p < 1$, and use $\mathrm{X}^F = \{x \in \mathcal{X}|F_\phi(x) < p\}$ as an approximate feasible region. A case where this approximation is valid is that constraint violations can be avoided in an infinite horizon as long as they can be avoided in a finite number of steps [40]. In our experiments, we found that a constant value $p = 0.1$ can achieve zero constraint violation on all classic control tasks. Inside the feasible region, we update the policy according to Eq. (16). Since the policy from the previous iteration is feasible inside its own feasible region, it can be viewed as a feasible initial point. This allows us to use the interior point method to solve the optimization problem in Eq. (16). We apply a logarithm function to the constraint and add it to the objective function, obtaining the following policy loss:

$$L_{\pi,\mathrm{f}}(\theta) = \mathbb{E}_{x\sim\mathcal{D},u\sim\pi_\theta(\cdot|x)}\{\mathbf{1}[G_\phi(x, u) < p](\alpha\log\pi_\theta(u|x)-$$
$$\min_{i\in\{1,2\}} Q_{\omega_i}(x, u) - 1/t \cdot \log(p - G_\phi(x, u)))\}. \quad (34)$$

We increase the weight $t$ by a factor every several iterations, which is a standard practice in the interior point method. As $t \to \infty$, the minima of Eq. (34) approaches the solution of Eq. (16). Compared with other constrained optimization methods such as the dual ascent, the advantage of the interior point method is that its intermediate solutions are always feasible regardless of the value of $t$. This ensures that the monotonic expansion property of the feasible region is preserved. Outside the feasible region, we minimize the action feasibility function:

$$L_{\pi,\mathrm{i}}(\theta) = \mathbb{E}_{x\sim\mathcal{D},u\sim\pi_\theta(\cdot|x)}\{G_\phi(x, u)\}. \quad (35)$$

With (34) and (35), the total policy loss is

$$L_\pi(\theta) = L_{\pi,\mathrm{f}}(\theta) + L_{\pi,\mathrm{i}}(\theta). \quad (36)$$

The loss function of the temperature $\alpha$ is the same as that in SAC:

$$L(\alpha) = \mathbb{E}_{x \sim D} \left\{ -\alpha \log \pi_\theta(u|x) - \alpha \bar{\mathcal{H}} \right\}, \qquad (37)$$

where $\bar{\mathcal{H}}$ is the target entropy.

## VII. Experiments

We aim to answer the following two questions through our experiments:

1) Does FPI outperform existing safe RL algorithms in terms of safety and optimality?
2) Does FPI achieve monotonic feasible region expansion and monotonic value function improvement, and eventually converge to the maximum feasible region and optimal value function?

To answer Question 1, we perform extensive experiments on twelve safe RL environments, including grid world tasks, classic control tasks, and high-dimensional robot navigation tasks. To answer Question 2, we visualize the feasible region and state-value function during training process to check monotonicity, and we also compare the final feasible regions with the maximum ones and the final reward performance with that of an optimal controller to check optimality.

### A. Environments

We choose four grid world tasks, four classic control tasks, and four high-dimensional robot navigation tasks. The grid world tasks are borrowed from AI Safety Gridworlds [41] with slight modifications, as shown in Fig. 3. In **SafeInterruptibility**, the agent should first press the button to disable the interruption, which is defined as a constraint, and then pass through to reach the goal. In **SideEffects**, the agent should push the box to clear a path to the goal. The constraint is that the box should not be pushed to a corner. In **Supervisor**, the agent receives a negative reward for stepping on the punishment with a 50% probability. The constraint requires the agent to always avoid the punishment. In **BoatRace**, the agent is rewarded for moving clockwise through the checkpoints. However, revisiting the same checkpoint back and forth results in constraint violations. For more details about the tasks, please refer to the original paper [41].
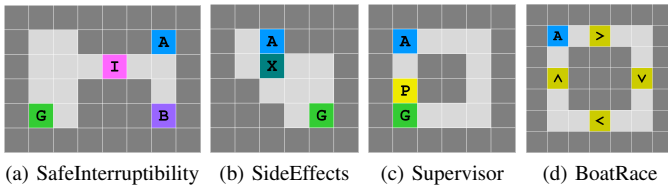


Fig. 3. Four grid world tasks. **A** denotes the agent. **G** denotes the goal. **I** in (a) denotes the interruption. **B** in (a) denotes the button. **X** in (b) denotes the box. **P** in (c) denotes the punishment. The arrows in (d) denote the checkpoints.

The classic control tasks include regulation and trajectory tracking, covering both linear and nonlinear systems, as shown in Fig. 4. **ACC** requires controlling a vehicle to follow a leading vehicle while keeping the distance within given limits. **LK** requires keeping a vehicle in a lane by controlling its steering angle. **Pendulum** requires swinging a pendulum to the upright position while keeping its angular position within given limits. **Quadrotor** is borrowed from safe-control-gym [42], where a 2D quadrotor is required to follow a circular trajectory in the vertical plane while keeping the vertical position with given limits. We defer the details of these environments, including the design of their state spaces, action spaces, dynamics, reward functions, and constraints, to Appendix I.
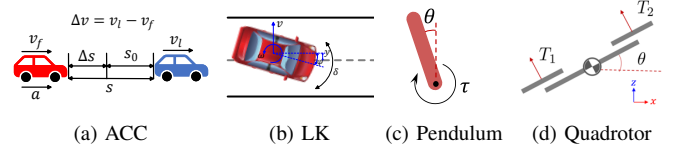


Fig. 4. Schematics of four classic control tasks.

The robot navigation tasks are borrowed from Safety Gym [38], as shown in Fig. 5. **PointGoal** and **CarGoal** require the agent to reach the goal while avoiding hazards, and **PointPush** and **CarPush** require the agent to push the box to the goal.
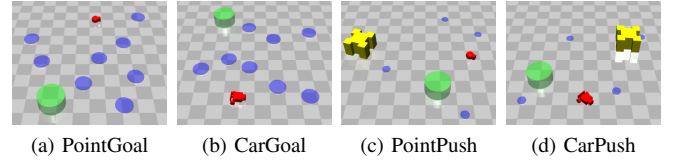


Fig. 5. Snapshots of four robot navigation tasks. The robots are in red. The blue circles represent the hazards. The Green cylinders represent the goals. The yellow objects represent the boxes.

### B. Results

For grid world tasks, we adopt the framework of Q-learning [43], which uses a table to store the values of all state-action pairs. Apart from FPI, which in this case we denote as **Q-learning-FPI**, we include another two algorithms for comparison: vanilla **Q-learning** and Q-learning with the Lagrange method, which we denote as **Q-learning-Lag**. Both Q-learning-FPI and Q-learning-Lag maintain another Q table for storing cost values. Q-learning-Lag selects the optimal action from a weighted sum of two Q tables, with the weight adjusted by dual ascent. Q-learning-FPI uses the cost Q table to eliminate infeasible actions and selects the optimal one from the rest using the reward Q table. Hyperparameters can be found in Appendix III.

For classic control tasks and robot navigation tasks, we compare FPI with (1) iterative unconstrained RL algorithms: SAC Lagrangian (**SAC-Lag**) [44] and SAC with CBF as constraint (**SAC-CBF**), and (2) constrained policy optimization methods: **CPO** [12] and PPO Lagrangian (**PPO-Lag**) [38]. The CBFs are handcrafted and described in Appendix II.

*1) Training curves:* The training curves of average episode cost and average episode return on grid world tasks are shown in Fig. 6. The superiority of FPI over the Lagrange

method is evident: Q-learning-FPI quickly converges to zero constraint violation and optimal return across all tasks, while Q-learning-Lag exhibits slower convergence and significant oscillations. The reason is that FPI immediately eliminates infeasible actions once they are identified by the feasibility function. In contrast, the Lagrange method requires iteratively updating the multiplier until the cost Q value overweighs the reward Q value to eliminate infeasible actions.

The training curves on classic control and robot navigation tasks are shown in Fig. 7. FPI-SAC achieves comparable or higher performance than the baselines while keeping near-zero constraint violations on all tasks. SAC-CBF also performs fairly safe, but this comes at the cost of sacrificing optimality. In comparison, the episode cost curves of other algorithms either fail to converge to zero or fluctuate severely. For example, CPO acts highly unsafely on LK and all robot navigation tasks except CarPush, so does SAC-Lag on two Point robot tasks, and PPO-Lag on LK and Quadrotor. These are due to the impacts of CPO's infeasibility problem and the Lagrange multiplier method's lack of monotonicity guarantee.

*2) Comparison with optimal policy:* For four classic control tasks, we use model predictive control (MPC) with a long enough horizon to approximate an optimal controller and compare the five algorithms with the MPC controller. Evaluation metrics include constraint violation ratio $R_{\text{vio}} = N_{\text{vio}}/N$, and average normalized return $R_{\text{norm}} = (R_{\text{alg}} - R_{\text{base}})/(R_{\text{MPC}} - R_{\text{base}}) - 1$, where $N_{\text{vio}}$ is the number of episodes where the constraint is violated, and $N$ is the total number of episodes. We uniformly discretize the state space to form a grid of feasible initial states, each point of which yields an episode. $R_{\text{alg}}$, $R_{\text{base}}$ and $R_{\text{MPC}}$ are the average return across $N$ episodes of a certain algorithm, a random policy and an MPC controller, respectively. The results are averaged across three random seeds, as listed in Table I. We further apply an MPC controller without considering any constraint to show that the constraints do take effect. Results show that FPI-SAC achieves zero constraint violation on all four tasks, while all other algorithms violate constraints on one or more tasks. Meanwhile, FPI-SAC ranks high in terms of return on all tasks.

*3) Visualization:* To show the monotonic expansion of the feasible region and monotonic improvement of the state-value function, we visualize some intermediate training results of FPI-SAC, as shown in Fig. 8. The state values are presented as heat maps. The maximum feasible regions and the current feasible regions are outlined with black lines and green lines, respectively. In all tasks, FPI-SAC demonstrates monotonicity in terms of both the feasible region and state value, and the feasible regions gradually converge to the maximum ones. Note that the initial state values in Fig. 8(b) and 8(c) are high only because the value networks are initialized to produce near-zero outputs and haven't acquired to give the true values at this early stage.

## VIII. CONCLUSION

This paper proposes FPI, the first foundational dynamic programming algorithm for safe RL. FPI follows the ACS iteration scheme and iteratively performs PEV, RID, and region-wise PIM. We prove that during the process of FPI, the feasible region monotonically expands, the state-value function monotonically increases within the feasible region, and both of them geometrically converge to their optima. Based on FPI, we develop an algorithm with neural network approximation to address continuous state and action spaces. Experiments show that our algorithm learns strictly safe and near-optimal policies with accurate feasible regions on low-dimensional tasks, and achieves lower constraint violations and better performance than existing algorithms on high-dimensional tasks.

## APPENDIX I
### DETAILS OF CLASSIC CONTROL TASKS

In ACC, the system dynamics is

$$\begin{pmatrix} \dot{x}_1 \\ \dot{x}_2 \end{pmatrix} = \begin{pmatrix} x_2 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ -1 \end{pmatrix} \mathbf{u},$$

where $\mathbf{x} = [x_1, x_2]^\top \triangleq [\Delta s, \Delta v]^\top$, with $\Delta s = s - s_0$ standing for the difference between actual distance $s$ and expected distance $s_0$ between the two vehicles and $\Delta v$ standing for the relative velocity. The action $\mathbf{u} \triangleq [a]$ is the acceleration of the following vehicle. The reward function is defined as $r(\mathbf{x}, \mathbf{u}) = -0.001\Delta s^2 - 0.01\Delta v^2 - a^2$, and the constraint function is $h(\mathbf{x}) = |\Delta s| - \Delta s_{\max}$.

In LK, the state $\mathbf{x} \triangleq [y, \varphi, v, \omega]^\top$ includes $y$ the vertical position, $\varphi$ the yaw angle, $v$ the lateral velocity and $\omega$ the yaw rate. The vehicle follows a 2DoF bicycle model with the action be its steering angle, i.e. $\mathbf{u} \triangleq [\delta]$. The reward function penalizes stabilization errors and aggressive behaviour: $r(\mathbf{x}, \mathbf{u}) = -0.01y^2 - 0.01\varphi^2 - v^2 - \omega^2 - \delta^2$, and the constraint function is $h(\mathbf{x}) = |y| - L/2$, where $L$ is the width of the lane.

Pendulum is borrowed from the RL benchmark Gymnasium [45] with an additional constraint $h(\mathbf{x}) = |\theta| - \theta_{\max} \leq 0$.

In Quadrotor, the state is $\mathbf{x} \triangleq [x, \dot{x}, z, \dot{z}, \theta, \dot{\theta}]^\top$, where $(x, z)$ is the position of the quadrotor on $xz$-plane, and $\theta$ is the pitch angle. The action of the system is $\mathbf{u} \triangleq [T_1, T_2]^\top$, including the thrusts generated by two pairs of motors. The reward is:

$$r(\mathbf{x}, \mathbf{u}) = -\|[x - x_{\text{ref}}, z - z_{\text{ref}}]\|_2^2 - 0.1\theta^2 - \dot{\theta}^2 \\ - 0.1(T_1 - T_0)^2 - 0.1(T_2 - T_0)^2,$$

where $(x_{\text{ref}}, z_{\text{ref}})$ is the reference position the quadrotor is supposed to be at, and $T_0$ is the thrust needed for balancing the gravity. The reference position moves along a circle $x^2 + z^2 = 0.5^2$ with a constant angular velocity. The constraint function $h(\mathbf{x}) = \max\{|z| - z_{\max}, |\theta| - \theta_{\max}, \text{err} - \text{err}_{\max}\}$, where $\text{err} = \|[x - x_{\text{ref}}, z - z_{\text{ref}}]\|_2$, restricts the quadrotor to stay in a rectangular area.

## APPENDIX II
### HANDCRAFTED CBFs

With a CBF $B(\mathbf{x}) : \mathcal{X} \to \mathbb{R}$, the constraint is constructed as $B(f(\mathbf{x}, \mathbf{u})) - B(\mathbf{x}) \leq -\lambda B(\mathbf{x})$, where $\lambda > 0$ is a constant. The CBF for ACC is

$$B(\mathbf{x}) = \begin{cases} -10 + \Delta s + 4.5\Delta v & \Delta v \geq 0 \\ -10 - \Delta s - 3.2\Delta v & \Delta v < 0 \end{cases}.$$
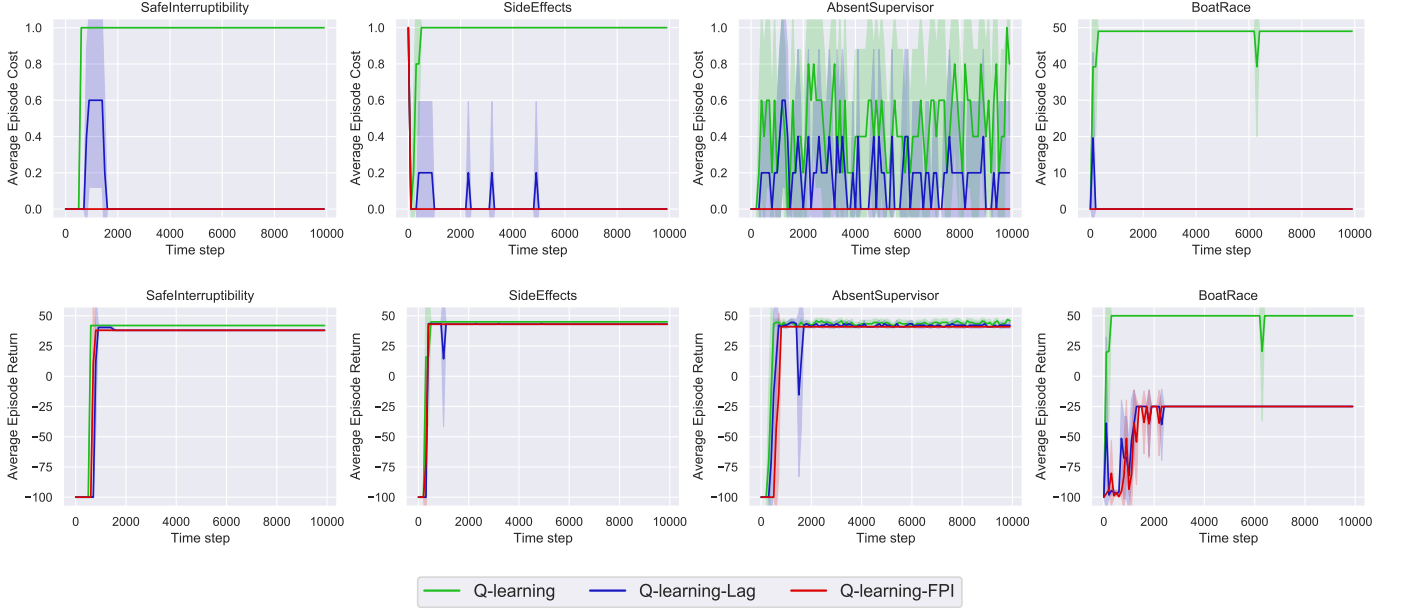
Fig. 6. Training curves on grid world tasks. The solid lines represent the mean and the shaded areas represent the 95% confidence intervals over 5 seeds. The first row shows the average episode cost, and the second row shows the average episode return.

TABLE I
CONSTRAINT VIOLATION AND NORMALIZED RETURN ON FOUR CLASSIC CONTROL TASKS.

| | ACC | | LK | | Pendulum | | Quadrotor | |
|---|---|---|---|---|---|---|---|---|
| | $R_{\text{vio}}$ | $R_{\text{norm}}$ | $R_{\text{vio}}$ | $R_{\text{norm}}$ | $R_{\text{vio}}$ | $R_{\text{norm}}$ | $R_{\text{vio}}$ | $R_{\text{norm}}$ |
| CPO | $4.17\pm0.00$ | $0.17\pm0.14$ | $100.00\pm0.00$ | $-17.56\pm1.59$ | $0.00\pm0.00$ | $-0.04\pm0.01$ | $0.00\pm0.00$ | $-0.06\pm0.01$ |
| SAC-Lag | $0.00\pm0.00$ | $-0.20\pm0.17$ | $0.00\pm0.00$ | $-0.22\pm0.06$ | $0.00\pm0.00$ | $-0.59\pm0.81$ | $8.89\pm3.85$ | $-5.81\pm1.55$ |
| PPO-Lag | $2.78\pm2.41$ | $-1.54\pm0.44$ | $8.70\pm0.00$ | $-0.36\pm0.60$ | $2.90\pm2.51$ | $-0.62\pm0.22$ | $100.00\pm0.00$ | $-94.34\pm1.37$ |
| SAC-CBF | $4.17\pm0.00$ | $-0.07\pm0.05$ | $0.00\pm0.00$ | $-3.84\pm0.34$ | $8.70\pm0.00$ | $-1.34\pm1.15$ | $31.11\pm6.29$ | $-28.43\pm5.37$ |
| FPI-SAC (ours) | $0.00\pm0.00$ | $-0.01\pm0.01$ | $0.00\pm0.00$ | $-0.25\pm0.17$ | $0.00\pm0.00$ | $-0.10\pm0.02$ | $0.00\pm0.00$ | $-0.03\pm0.01$ |
| MPC w/o cstr. | $20.83\pm0.00$ | $0.49\pm0.00$ | $43.48\pm0.00$ | $0.17\pm0.00$ | $17.39\pm0.00$ | $1.14\pm0.00$ | $100.00\pm0.00$ | $0.03\pm0.00$ |

[1] The values denote mean±standard deviation.
[2] MPC w/o cstr. refers to an MPC controller only maximizing return without considering constraints.

The CBF for LK is

$$B(\mathbf{x}) = \begin{cases} -L/2 + 1.8y + 3.2\varphi^2 + 0.5v + 0.8\omega & \phi \geq 0 \\ -L/2 - 1.8y - 3.2\varphi^2 - 0.5v - 0.8\omega & \phi < 0 \end{cases}.$$

The CBF for Pendulum is

$$B(\mathbf{x}) = \begin{cases} -\theta_{\max} + \theta + 0.3\dot{\theta} & \dot{\theta} \geq 0 \\ -\theta_{\max} - \theta - 0.3\dot{\theta} & \dot{\theta} < 0 \end{cases}.$$

The CBF for the Quadrotor is

$$B(\mathbf{x}) = -z_{\max} + 2.35z^2$$
$$+ \begin{cases} 2\text{err}^2 + 0.2\dot{z} & 0 \leq z \leq z_{\max} \\ 2\text{err}^2 - 0.2\dot{z} & -z_{\max} \leq z < 0 \\ 2\text{err}^2 & \text{others} \end{cases}.$$

The CBF for Safety Gym environments is $B(\mathbf{x}) = r - d - 0.3\dot{d}$, where $r$ is the radius of hazards and $d$ is the distance to the center of the nearest hazard.

TABLE II
HYPERPARAMETERS FOR GRID WORLD TASKS

| Hyperparameter | Value |
|---|---|
| Discount factor | 0.9 |
| Learning rate | 0.5 |
| Random exploration probability ($\epsilon$-greedy) | 0.1 |

APPENDIX III
HYPERPARAMETERS

REFERENCES

[1] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev *et al.*, "Grandmaster level in starcraft ii using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.

[2] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver, "Mastering atari, go, chess and shogi by planning with a learned model," *Nature*, vol. 588, no. 7839, pp. 604–609, 2020.

[3] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray *et al.*, "Learning dexterous in-hand manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 3–20, 2020.
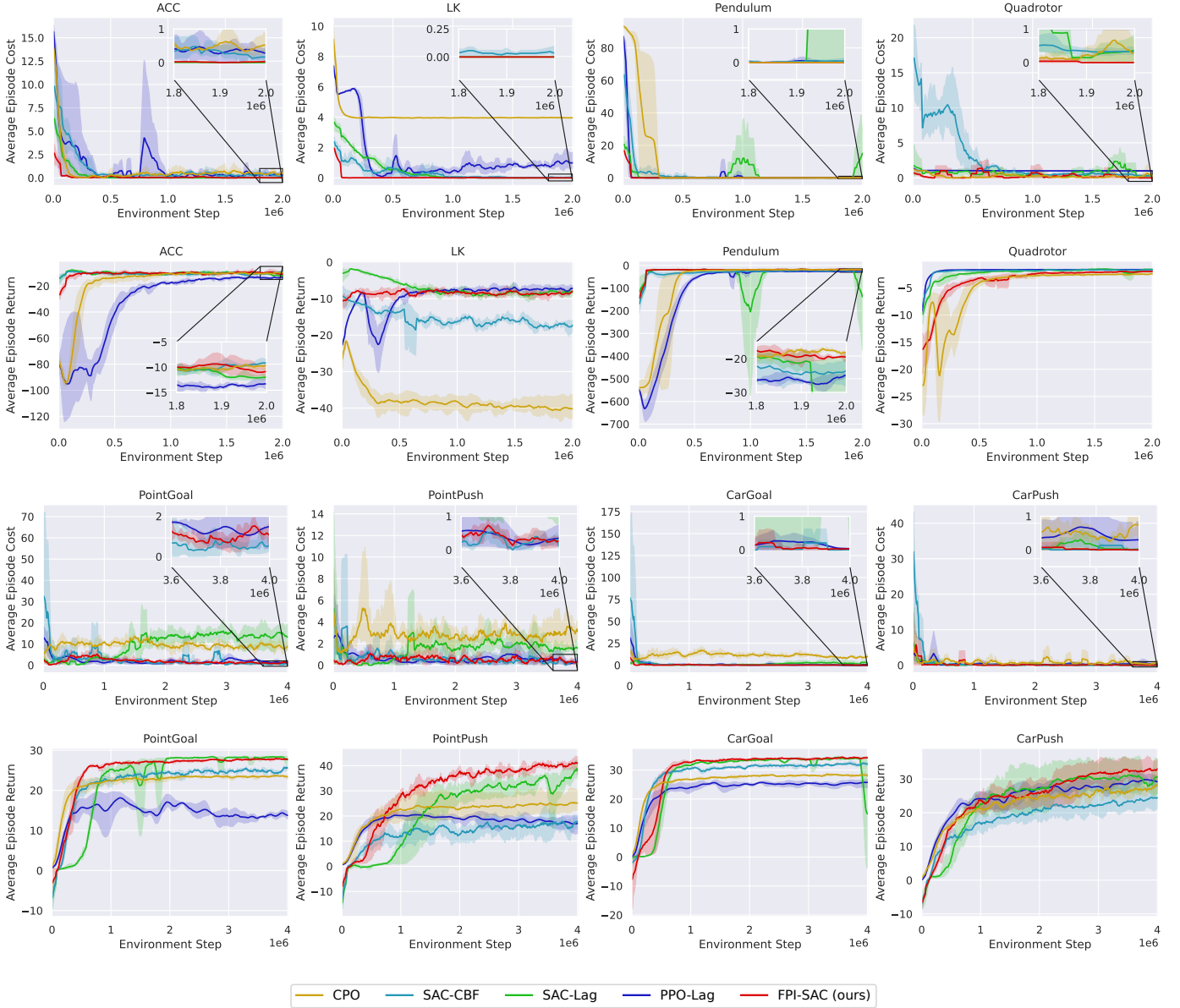
Fig. 7. Training curves on classic control and robot navigation tasks. The shaded areas represent the 95% confidence intervals over 3 seeds.

TABLE III
HYPERPARAMETERS FOR CONTINUOUS SPACE TASKS

| Hyperparameter | Value | |
| --- | --- | --- |
| | Classic | Safety Gym |
| Discount factor | 0.99 | |
| Number of hidden layers | 2 | |
| Number of hidden neurons | 256 | |
| Optimizer | Adam($\beta_1$=0.99,$\beta_2$=0.999) | |
| Activation function | ReLU | |
| Target entropy | -dim($\mathcal{U}$) | |
| Initial temperature | 1.0 | |
| Target smoothing coefficient | 0.005 | |
| Learning rate | 1e-4 | |
| Batch size | 256 | 1024 |
| Replay buffer size | $2 \times 10^6$ | $4 \times 10^6$ |
| Feasibility threshold ($p$) | 0.1 | |
| Initial $t$ | 1.0 | |
| $t$ increase factor | 1.1 | |
| $t$ update delay | 10000 | |

[4] Y. Guan, Y. Ren, Q. Sun, S. E. Li, H. Ma, J. Duan, Y. Dai, and B. Cheng, "Integrated decision and control: Toward interpretable and computationally efficient driving intelligence," *IEEE Transactions on Cybernetics*, pp. 1–15, 2022.

[5] S. E. Li, *Reinforcement learning for sequential decision and optimal control*. Springer, 2023.

[6] R. S. Sutton, "Reinforcement learning: An introduction," *A Bradford Book*, 2018.

[7] T. Lillicrap, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.

[8] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *International conference on machine learning*. PMLR, 2018, pp. 1587–1596.

[9] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.

[10] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

[11] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International Conference on Machine*
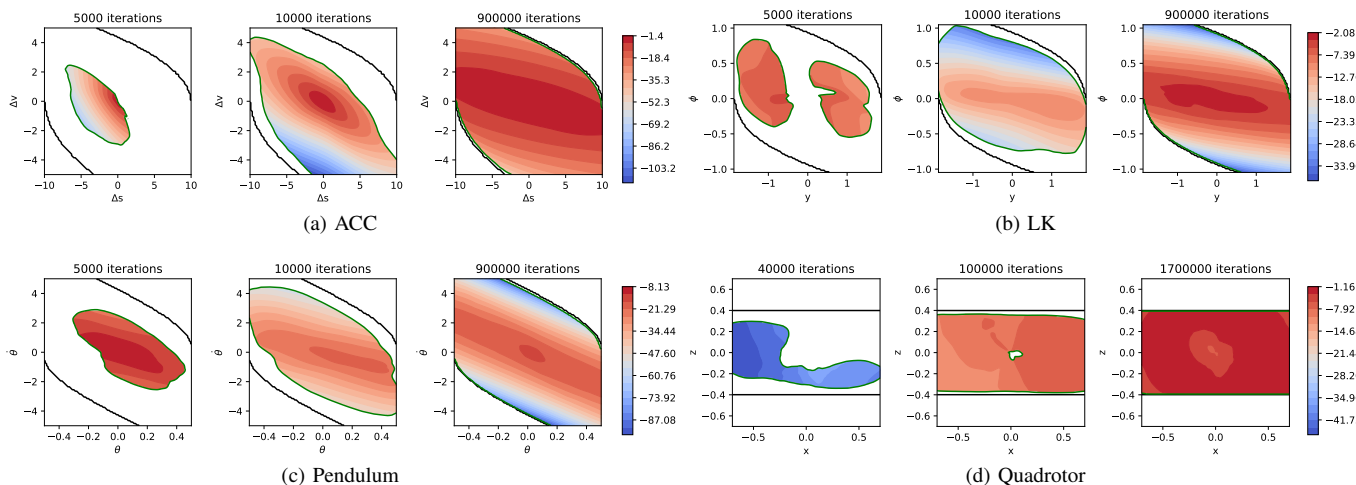
Fig. 8. Heat maps and zero-sublevel sets of the feasibility functions on four classic control tasks. The heat maps correspond to the state values. The black lines are boundaries of the maximum feasible regions, while the green ones around the colored regions are boundaries of zero contours of the feasibility functions.

*Learning*. PMLR, 2015, pp. 1889–1897.

[12] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," in *International Conference on Machine Learning*. PMLR, 2017, pp. 22–31.

[13] S. Paternain, L. Chamon, M. Calvo-Fullana, and A. Ribeiro, "Constrained reinforcement learning has zero duality gap," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[14] Y. Chow, M. Ghavamzadeh, L. Janson, and M. Pavone, "Risk-constrained reinforcement learning with percentile risk criteria," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6070–6120, 2017.

[15] C. Tessler, D. J. Mankowitz, and S. Mannor, "Reward constrained policy optimization," in *International Conference on Learning Representations*, 2019.

[16] D. Yu, W. Zou, Y. Yang, H. Ma, S. E. Li, Y. Yin, J. Chen, and J. Duan, "Safe model-based reinforcement learning with an uncertainty-aware reachability certificate," *IEEE Transactions on Automation Science and Engineering*, 2023.

[17] J. F. Fisac, N. F. Lugovoy, V. Rubies-Royo, S. Ghosh, and C. J. Tomlin, "Bridging hamilton-jacobi safety analysis and reinforcement learning," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 8550–8556.

[18] K. C. Hsu, V. Rubies-Royo, C. J. Tomlin, and J. F. Fisac, "Safety and liveness guarantees through reach-avoid reinforcement learning," in *17th Robotics: Science and Systems, RSS 2021*. MIT Press Journals, 2021.

[19] D. Yu, H. Ma, S. Li, and J. Chen, "Reachability constrained reinforcement learning," in *International Conference on Machine Learning*. PMLR, 2022, pp. 25 636–25 655.

[20] M. Ganai, Z. Gong, C. Yu, S. Herbert, and S. Gao, "Iterative reachability estimation for safe reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[21] H. Ma, J. Chen, S. Eben, Z. Lin, Y. Guan, Y. Ren, and S. Zheng, "Model-based constrained reinforcement learning using generalized control barrier function," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 4552–4559.

[22] J. Wang, D. Zhang, J. Zhang, H. Zhu, S. Hu, and C. Qin, "Safe adaptive dynamic programming method for nonlinear safety-critical systems with disturbance," in *2021 6th International Conference on Robotics and Automation Engineering (ICRAE)*, 2021, pp. 97–101.

[23] J. Xu, J. Wang, J. Rao, Y. Zhong, and H. Wang, "Adaptive dynamic programming for optimal control of discrete-time nonlinear system with state constraints based on control barrier function," *International Journal of Robust and Nonlinear Control*, vol. 32, no. 6, pp. 3408–3424, 2022.

[24] Y. Yang, Y. Jiang, Y. Liu, J. Chen, and S. E. Li, "Model-free safe reinforcement learning through neural barrier certificate," *IEEE Robotics and Automation Letters*, 2023.

[25] C. Mu, K. Wang, X. Xu, and C. Sun, "Safe adaptive dynamic programming for multiplayer systems with static and moving no-entry regions," *IEEE Transactions on Artificial Intelligence*, vol. 5, no. 5, pp. 2079–2092, 2024.

[26] Y. Yang, Y. Zhang, W. Zou, J. Chen, Y. Yin, and S. Eben Li, "Synthesizing control barrier functions with feasible region iteration for safe reinforcement learning," *IEEE Transactions on Automatic Control*, vol. 69, no. 4, pp. 2713–2720, 2024.

[27] C. Liu and M. Tomizuka, "Control in a safe set: Addressing safety in human-robot interactions," in *Dynamic Systems and Control Conference*, vol. 46209. American Society of Mechanical Engineers, 2014, p. V003T42A003.

[28] H. Ma, C. Liu, S. E. Li, S. Zheng, and J. Chen, "Joint synthesis of safety certificate and safe control policy using constrained reinforcement learning," in *The 4th Annual Learning for Dynamics and Control Conference*, vol. 168. PMLR, 2022, pp. 97–109.

[29] S. M. Richards, F. Berkenkamp, and A. Krause, "The lyapunov neural network: Adaptive stability certification for safe learning of dynamical systems," in *Conference on Robot Learning*. PMLR, 2018, pp. 466–476.

[30] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh, "A lyapunov-based approach to safe reinforcement learning," *Advances in neural information processing systems*, vol. 31, 2018.

[31] Y.-C. Chang, N. Roohi, and S. Gao, "Neural lyapunov control," *Advances in neural information processing systems*, vol. 32, 2019.

[32] A. Stooke, J. Achiam, and P. Abbeel, "Responsive safety in reinforcement learning by pid lagrangian methods," in *International Conference on Machine Learning*. PMLR, 2020, pp. 9133–9143.

[33] B. Peng, J. Duan, J. Chen, S. E. Li, G. Xie, C. Zhang, Y. Guan, Y. Mu, and E. Sun, "Model-based chance-constrained reinforcement learning via separated proportional-integral lagrangian," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[34] T.-Y. Yang, J. Rosca, K. Narasimhan, and P. J. Ramadge, "Projection-based constrained policy optimization," in *International Conference on Learning Representations*, 2020.

[35] Y. Zhang, Q. Vuong, and K. Ross, "First order constrained optimization in policy space," *Advances in Neural Information Processing Systems*, vol. 33, pp. 15 338–15 349, 2020.

[36] L. Yang, J. Ji, J. Dai, L. Zhang, B. Zhou, P. Li, Y. Yang, and G. Pan, "Constrained update projection approach to safe policy optimization," *Advances in Neural Information Processing Systems*, vol. 35, pp. 9111–9124, 2022.

[37] Y. Liu, J. Ding, and X. Liu, "Ipo: Interior-point policy optimization under constraints," *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 04, pp. 4940–4947, 2020.

[38] A. Ray, J. Achiam, and D. Amodei, "Benchmarking safe exploration in deep reinforcement learning," *arXiv preprint arXiv:1910.01708*, vol. 7, p. 1, 2019.

[39] J. Dai, J. Ji, L. Yang, Q. Zheng, and G. Pan, "Augmented proximal policy optimization for safe reinforcement learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 6, pp. 7288–7295, 2023.

[40] G. Thomas, Y. Luo, and T. Ma, "Safe reinforcement learning by

imagining the near future," *Advances in Neural Information Processing Systems*, vol. 34, pp. 13 859–13 869, 2021.

[41] J. Leike, M. Martic, V. Krakovna, P. A. Ortega, T. Everitt, A. Lefrancq, L. Orseau, and S. Legg, "Ai safety gridworlds," *arXiv preprint arXiv:1711.09883*, 2017.

[42] Z. Yuan, A. W. Hall, S. Zhou, L. Brunke, M. Greeff, J. Panerati, and A. P. Schoellig, "Safe-control-gym: A unified benchmark suite for safe learning-based control and reinforcement learning in robotics," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 11 142–11 149, 2022.

[43] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, pp. 279–292, 1992.

[44] S. Ha, P. Xu, Z. Tan, S. Levine, and J. Tan, "Learning to walk in the real world with minimal human effort," in *Conference on Robot Learning*. PMLR, 2021, pp. 1110–1120.

[45] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.