

Joint Synthesis of Safety Certificate and Safe Control Policy using Constrained Reinforcement Learning

Haitong Ma

MAHT19@MAILS.TSINGHUA.EDU.CN

School of Vehicle and Mobility, Tsinghua University, Beijing 100084, China

Changliu Liu

CLIU6@ANDREW.CMU.EDU

Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA

Shengbo Eben Li

LISHBO@TSINGHUA.EDU.CN

Sifa Zheng

ZSF@TSINGHUA.EDU.CN

School of Vehicle and Mobility, Tsinghua University, Beijing 100084, China

Jianyu Chen

JIANYUCHEN@TSINGHUA.EDU.CN

Institute of Interdisciplinary Information Science, Tsinghua University, Beijing 100084, China

Shanghai Qizhi Institute, Shanghai 200000, China

Editors: R. Firoozi, N. Mehr, E. Yel, R. Antonova, J. Bohg, M. Schwager, M. Kochenderfer

Abstract

Safety is the major consideration in controlling complex dynamical systems using reinforcement learning (RL), where the *safety certificates* can provide provable safety guarantees. A valid safety certificate is an energy function indicating that safe states are with low energy, and there exists a corresponding *safe control policy* that allows the energy function to always dissipate. The safety certificates and the safe control policies are closely related to each other and both challenging to synthesize. Therefore, existing learning-based studies treat either of them as prior knowledge to learn the other, limiting their applicability to general systems with unknown dynamics. This paper proposes a novel approach that simultaneously synthesizes the energy-function-based safety certificates and learns the safe control policies with constrained reinforcement learning (CRL). We do not rely on prior knowledge about *either* a prior control law *or* a perfect safety certificate. In particular, we formulate a loss function to optimize the safety certificate parameters by minimizing the occurrence of energy increases. By adding this optimization procedure as an outer loop to the Lagrangian-based CRL, we jointly update the policy and safety certificate parameters, and prove that they will converge to their respective local optima, the optimal safe policies and valid safety certificates. Finally, we evaluate our algorithms on multiple safety-critical benchmark environments. The results show that the proposed algorithm learns solidly safe policies with no constraint violation. The validity or feasibility of synthesized safety certificates is also verified numerically.

Keywords: Safety Certificates, Safe Control, Energy functions, Safety Index Synthesis, Constrained Reinforcement Learning, Multi-Objective Learning

1. Introduction

Safety is critical when applying state-of-the-art artificial intelligence studies to real-world applications, like autonomous driving (Sallab et al., 2017; Chen et al., 2021), robotic control (Richter et al., 2019; Ray et al., 2019). Safe control is one of the most common tasks among these real-world applications, requiring that the hard safety constraints must be obeyed persistently. However, learning

a safe control policy is hard for the naive trial-and-error mechanism of RL since it penalizes the dangerous actions *after* experiencing them.

Meanwhile, in the control theory, there exist studies about *energy-function-based* provable safety guarantee of dynamic systems called the *safety certificate*, or *safety index* (Wieland and Allgöwer, 2007; Ames et al., 2014; Liu and Tomizuka, 2014). These methods first synthesize energy functions such that the safe states have low energy, then design control laws satisfying the *safe action constraints* to make the systems dissipate energy (Wei and Liu, 2019). If there exists a feasible policy for *all states* in a safe set to satisfy the *safe action constraints* dissipating the energy, the system will never leave the safe set (i.e., forward invariance). Despite its soundness, the safety index synthesis (SIS) by hand is extremely hard for complex systems, which stimulates a rapidly growing interest in learning-based SIS (Chang et al., 2020; Saveriano and Lee, 2019; Srinivasan et al., 2020; Ma et al., 2021a; Qin et al., 2021). Nevertheless, These studies usually require known dynamical models (white-box, black-box or learning-calibrated) to design control laws. Furthermore, obtaining the policy satisfying safe action constraints is also challenging. Adding *safety shields or layers* to obtain supervised RL policies is a common approach (Wang et al., 2017; Agrawal and Sreenath, 2017; Cheng et al., 2019; Taylor et al., 2020), but these studies usually assume to know the valid safety certificates.

In general safe control tasks with unknown dynamics, one usually has access to *neither* the control laws *nor* perfect safety certificates, which makes the previous two kinds of studies fall into a paradox—they rely on each other as the prior knowledge. Therefore, this paper proposes a novel algorithm without prior knowledge about model-based control laws or valid safety certificates. We define a loss function for SIS by minimizing the occurrence of energy increases. Then we formulate a CRL problem (rather than the commonly used shield methods) to unify the loss functions of SIS and CRL. By adding SIS as an outer loop to the Lagrangian-based solution to CRL, we jointly update the policies and safety certificates, and prove that they will converge to their respective local optima, the optimal safe policies and the valid safety certificates.

Contributions. Our main contributions are: 1. We propose an algorithm of joint CRL and SIS that learns the safe policies and synthesizes the safety certificates simultaneously. This is the first algorithm requiring no prior knowledge of control laws or valid safety certificates. 2. We unify the loss function formulations of SIS and CRL. We therefore can form the multi-timescale adversarial RL training and prove its convergence. 3. We evaluate the proposed algorithm on multiple safety-critical benchmark environments Results demonstrate that we can simultaneously synthesize valid safety certificates and learn safe policies with zero constraint violation.

1.1. Related works

Representative energy-function-based safety certificates include barrier certificates (Prajna et al., 2007), control barrier functions (CBF) (Wieland and Allgöwer, 2007), safety set algorithm (SSA) (Liu and Tomizuka, 2014) and sliding mode methods (Gracia et al., 2013). Recent learning-based studies can be mainly divided into *learning-based SIS* and *learning safe control policies supervised by certificates*. Chang et al. (2020); Luo and Ma (2021) use explicit models to rollout or project actions to satisfy safe action constraints. Jin et al. (2020); Qin et al. (2021) guide certificate learning with LQR controllers, Zhao et al. (2021) requires a black-box model to query online, and Saveriano and Lee (2019); Srinivasan et al. (2020) use labeled data to fit certificates with supervised learning. The latter one, learning safe policy with supervisory usually assumed a valid safety certificate (Wang

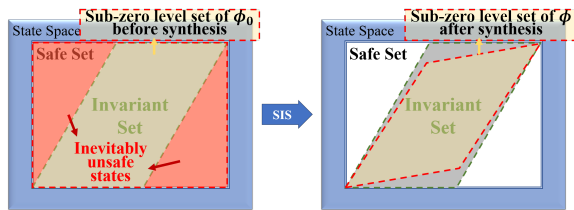


Figure 1: Safety index synthesis (SIS). Inevitably unsafe states should be excluded during the SIS.

et al., 2017; Cheng et al., 2019; Taylor et al., 2020). It’s a natural thought that one could learn the dynamic models to handle these issues (like Cheng et al., 2019; Luo and Ma, 2021), but learning models is much more complex than only learning policies, especially in RL tasks.

2. Problem Formulations

We consider the safety specification that the system state s should be constrained in a connected and closed set \mathcal{S}_s which is called the *safe set*. \mathcal{S}_s should be a zero-sublevel set of a safety index function $\phi_0(\cdot)$ denoted by $\mathcal{S}_s = \{s | \phi_0(s) \leq 0\}$. We study the Markov decision process (MDP) with deterministic dynamics (a reasonable assumption when dealing with safe control problems), defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{F}, r, c, \gamma, \phi)$, where \mathcal{S}, \mathcal{A} is the state and action space, $\mathcal{F} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ is the unknown system dynamics, $r, c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward and cost function, γ is the discounted factor, and $\phi : \mathcal{S} \rightarrow \mathbb{R}$ is the energy-function-based safety certificate, or called the *safety index*. A safe control law with respect to safety index ϕ should keep the system energy low, ($\phi \leq 0$) and dissipate the energy when the system is at high energy ($\phi > 0$). We use s' to represent the next state for simplicity. Then we can get the *safe action constraint*:

Definition 1 (Safe action constraint) For a given safety index ϕ , the safe action constraint is

$$\phi(s') < \max\{\phi(s) - \eta_D, 0\} \quad (1)$$

where η_D is a slack variable controlling the descent rate of safety index.

Definition 2 (Valid safety certificate) If there always exists an action $a \in \mathcal{A}$ satisfying (1) at s , or the safe action set $\mathcal{U}_s(s) = \{a | \phi(s') < \max\{\phi(s) - \eta_D, 0\}\}$ is always nonempty, we say the safety index ϕ is a **valid**, or **feasible** safety certificate.

A straightforward approach is to use the ϕ_0 as the safety certificate. However, these safe action constraints are possibly not satisfied with all the states in \mathcal{S}_s as shown in Figure 1. This problem is common in real-world tasks with actuator saturation and high relative-degree from safety specifications to control inputs (i.e., $\|\frac{\partial \phi_0}{\partial u}\| = 0$). For example, if ϕ_0 measures the distance between two autonomous vehicles, the collision may be inevitable because the relative speed is too high and brake force is limited. In this case, \mathcal{S}_s includes inevitably unsafe states. We need to assign high energy values to these inevitably-unsafe states, for example, by linearly combining the ϕ_0 and its high-order derivatives (Liu and Tomizuka, 2014). The valid safety certificate will guarantee safety by ensuring the *forward invariance* of a subset of \mathcal{S}_s .

Lemma 3 (Forward invariance (Liu and Tomizuka, 2014)) Define the zero-sublevel set of a valid safety index ϕ as $\mathcal{S}_s^\phi = \{s | \phi(s) \leq 0\}$. If ϕ is a valid safety certificate, then there exist policies to guarantee the forward invariance of $\mathcal{S}_s^\phi \cap \mathcal{S}_s$.

We therefore can formulate the CRL problem by adding the safe action constraints to RL optimization objective:

$$\max_{\pi} \mathbb{E}_{\tau \sim \pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \right\} = \mathbb{E}_s \{v^\pi(s)\} \quad \text{s.t. } \phi(s') - \max\{\phi(s) - \eta_D, 0\} < 0, \forall s \in \mathcal{D} \quad (2)$$

where $v^\pi(s)$ is the state-value function of s , $\mathcal{D} = \{s | f(s) > 0\}$ is the set of all possible states (f is the state distribution density).

Remark 4 (2) has *state-dependent* constraints; it can not be solved by previous model-free CRL since their constraint objectives are defined on the *expectation* over \mathcal{D} (Uchibe and Doya, 2007; Achiam et al., 2017; Chow et al., 2017; Tessler et al., 2018; Ray et al., 2019; Stooke et al., 2020).

We solve (2) based on our previous framework to solve state-dependent safety constraints in (Ma et al., 2021b) using the Lagrange multiplier networks $\lambda(s)$ in a Lagrangian-based approach. The Lagrange function is (Ma et al., 2021b)

$$\mathcal{L}'(\pi, \lambda) = \mathbb{E}_s \left\{ -v^\pi(s) + \lambda(s) (\phi(s') - \max\{\phi(s) - \eta_D, 0\}) \right\} \quad (3)$$

We can solve (2) by locating the saddle point of $\mathcal{L}'(\pi, \lambda)$, $\max_{\lambda} \min_{\pi} \mathcal{L}'(\pi, \lambda)$.

3. Joint Synthesis of Safety Certificate and Safe Control Policy

The key idea of this section is to unify the loss functions of CRL and SIS; we provide theoretical analyses of their equivalence.

3.1. Loss Function for Safety Index Synthesis

We construct the loss for optimizing a parameterized safety index by a measurement of the *violation of constraint* (1)

$$J(\phi) = \mathbb{E}_s \left\{ [\phi(s') - \max\{\phi(s) - \eta_D, 0\}]^+ \right\} \Big|_{\pi=\pi^*} \quad (4)$$

where $[\cdot]^+$ means projecting the values to the positive half-space $[0, +\infty)$, π^* is the optimal safe policy (also a feasible policy when ϕ is a valid safety index) of (2), and $\cdot|_{\pi=\pi^*}$ represents the agent takes $\pi^*(s)$ to reach s' . Ideally, if ϕ is a valid safety index, there always exists control to satisfy (1), and $J(\phi) = 0$. For those imperfect ϕ , the inequality constraint in (2) may not hold for all states in \mathcal{D} , so we can optimize the loss to get better ϕ .

The joint synthesis algorithm is tricky since we need to handle *two different optimization problems*, (3) and (4). Recent similar studies integrate two optimizations by weighted sum (Qin et al., 2021) or alternative update (Luo and Ma, 2021), but their methods are more like intuitive approaches and lack a solid theoretical basis.

3.2. Unified Loss Function for Joint Synthesis

Lemma 5 (Statewise complementary slackness condition (Ma et al., 2021b)) For the problem (2), if the safe action set is not empty at state s , the optimal multiplier and optimal policy λ^*, π^* satisfies

$$\lambda^*(s)\{\phi(s') - \max\{\phi(s) - \eta_D, 0\}|_{\pi^*}\} = 0 \quad (5)$$

If the safe action set is empty at state s , then $\lambda^*(s) \rightarrow \infty$.

The lemma comes from the Karush-Kuhn-Tucker (KKT) necessary conditions for the problem (2).

Consider the Lagrange function (3) with the additional variable ϕ to optimize,

$$\mathcal{L}'(\pi, \lambda, \phi) = \mathbb{E}_s \{-v^\pi(s) + \lambda(s)(\phi(s') - \max\{\phi(s) - \eta_D, 0\})\} \quad (6)$$

we have the following lemma for the relationship between the loss function of policy and certificate synthesis

Lemma 6 If λ is clipped into a compact set $[0, \lambda_{\max}]$, where $\lambda_{\max} > \max_{s \in \{s | \mathcal{U}_s(s) \neq \emptyset\}} \lambda^*(s)$. Then

$$\mathcal{L}'(\pi^*, \lambda^*, \phi) = \lambda_{\max} J(\phi) + \Delta \quad (7)$$

where Δ is a constant irrelevant with ϕ .

Proof See Appendix A.1¹. ■

Theorem 7 (Unified loss for joint synthesis) The optimal safety certificate parameters with optimal policy-multiplier tuple in (6) is also the optimal safety certificate parameters under loss in (4)

$$\arg \min J(\phi) = \arg \min \mathcal{L}'(\pi^*, \lambda^*, \phi) \quad (8)$$

Proof See the *envelope theorem* on parameterized constraints in Milgrom and Segal (2002); Rockafellar (2015). ■

Finally, we unify the loss function of updating three elements: policy π , multiplier λ , and safety index function ϕ . The optimization problem is formulated by a multi-timescale adversarial training:

$$\min_{\phi} \max_{\lambda} \min_{\pi} \mathcal{L}'(\pi, \lambda, \phi) \quad (9)$$

4. Practical Algorithm using Constrained Reinforcement Learning

In this section, we explain the practical algorithm and convergence analysis.

1. The full paper with Appendix can be found on <https://arxiv.org/abs/2111.07695>.

4.1. Algorithm Details

The Lagrangian-based solution to CRL with statewise safety constraint is compatible with any existing unconstrained RL baselines, and we select the off-policy maximum entropy RL framework like soft actor-critic (SAC, Haarnoja et al., 2018a). According to (3), we need to add two neural networks to learn a state-action value function for safety index model $Q_\phi(s, a)$ (to approximate $\phi(s') - \max\{\phi(s) - \eta_D, 0\}$) and the multipliers $\lambda(s)$. Because the soft policy evaluation, or the soft Q-learning part, is the same as Haarnoja et al. (2018a), we only demonstrate the soft policy iteration part of the algorithm in Algorithm 1. We name the proposed algorithm as FAC-SIS². We denote the parameters of the policy network, multiplier network, safety index as θ, ξ, ζ , respectively. We use $G_\#, \# \in \{\theta, \xi, \zeta\}$ to denote the gradients to update policy, multiplier and certificate. Detailed computations of the gradients can be found in Appendix B.1. In addition, we assign multiple delayed updates (similar to Fujimoto et al., 2018), $m_\pi < m_\lambda < m_\phi$, to stabilize the adversarial optimizations.

Remark 8 A general parameterization rule of safety index is to linearly combine ϕ_0 and its high-order derivatives (Liu and Tomizuka, 2014), $\phi = \phi_0 + k_1\dot{\phi}_0 + \dots + k_n\phi_0^{(n)}$; the parameters ξ is $[k_1, k_2, \dots, k_n]$. How many high-order derivatives are needed depends on the system relative-degree. For example, the relative-degree of position constraints with force inputs is 2. This information should be included in observations of MDP. Otherwise, the observation can not fully describe how dangerous the agent is with respect to the safety constraint.

Algorithm 1 Soft Policy Improvement in FAC-SIS

Input: Buffer \mathcal{D} with sampled data, policy parameters θ , multiplier parameters ξ , safety index parameters ζ .

- 1: **if** gradient steps $\bmod m_\pi = 0$ **then** $\theta \leftarrow \theta - \overline{\beta_\pi} G_\theta$
- 2: **if** gradient steps $\bmod m_\lambda = 0$ **then** $\xi \leftarrow \xi + \overline{\beta_\lambda} G_\xi$
- 3: **if** gradient steps $\bmod m_\phi = 0$ **then** $\zeta \leftarrow \zeta - \overline{\beta_\zeta} G_\zeta$

Output: w_1, w_2, θ, ξ .

4.2. Convergence Analysis

The convergence proof of a three timescale adversarial training of (9) mainly follows the multi-timescale convergence according to Theorem 2 in Chapter 6 in Borkar (2009) about multiple timescale convergence of multi-variable optimization. Some studies also adopted this procedure to explain the convergence of RL algorithms from the perspective of stochastic optimization (Bhatnagar et al., 2009; Bhatnagar and Lakshmanan, 2012), especially those with Lagrangian-based methods (Chow et al., 2017). We incorporate the recent study on *clipped stochastic gradient descent* to further improve the generalization of this convergence proof (Zhang et al., 2019). We first give some assumptions:

Assumption 1 (learning rates) The learning rate schedules, $\{\beta_\theta(k), \beta_\xi(k), \beta_\zeta(k)\}$, satisfy

$$\sum_k \beta_\theta(k) = \sum_k \beta_\xi(k) = \sum_k \beta_\zeta(k) = \infty.$$

$$\sum_k \beta_\theta(k)^2, \sum_k \beta_\xi(k)^2, \sum_k \beta_\zeta(k)^2 < \infty, \beta_\xi(k) = o(\beta_\theta(k)), \beta_\zeta(k) = o(\beta_\xi(k)).$$

2. FAC refers to the FAC algorithm in our prior work (Ma et al., 2021b).

This assumption also implies that the policy converges in the fastest timescale, then the multipliers, and finally the safety index parameters.

Proposition 9 (Clipped gradient descent) *The actual learning rate used in Algorithm 1 is*

$$\overline{\beta}_{\#} := \min \{ \beta_{\#}, \beta_{\#} / \|G_{\#}\| \} \quad (10)$$

where $\# \in \{w_1, w_2, \theta, \xi, \zeta\}$ is the parameters, and $G_{\#}$ is the corresponding gradients.

Assumption 2 *The state and action are sampled from compact sets, and all neural networks are $L_0 - L_1$ smooth.*

As we are to finish a safe control problem that the agent should be confined in safe sets, and the actuator has physical limits, the bounded assumption is reasonable. We use multi-layer perceptron with continuous differentiable activation functions in practical implementations (details can be found in Appendix D).

Theorem 10 *Under all the aforementioned assumptions, the sequence of policy, multiplier, and safety index parameters tuple $(\theta_k, \xi_k, \zeta_k)$ converge almost surely to a locally optimal safety index parameters and its corresponding locally optimal policy and multiplier $(\theta^*, \xi^*, \zeta^*)$ as k goes to infinity.*

Proof See Appendix B.2. ■

5. Experiments

In our experiments, we focus on the following questions:

1. How does the proposed algorithm compare with other constraint RL algorithms? Can it achieve a safe policy with zero constraint violation?
2. How does the learning-based safety index synthesis outperform the handcrafted safety index or the original safety index in the safety performance?
3. Does the synthesized safety index allow safe control in all states the agent experienced?

To demonstrate the effectiveness of the proposed online synthesis rules, we select the safe reinforcement learning benchmark environments Safety Gym (Ray et al., 2019) with different tasks and obstacles. We name a specific environment by $\{\text{Obstacle type}\}-\{\text{Obstacle size}\}-\{\text{Task}\}$. We select six environments with different tasks and constraint objectives, where four of them are demonstrated in Figure 2, and others are provided in Appendix E.1.

Remark 11 *Some previous studies use controllers specially designed for goal-reaching tasks (Jin et al., 2020; Qin et al., 2021), while our algorithm can handle arbitrary complex tasks like the Push tasks.*

In this section, we use a fine-tuned form of safety index in [Zhao et al. \(2021\)](#)

$$\phi(s) = \sigma + d_{\min}^n - d^n - k\dot{d} \quad (11)$$

where d is the distance between the agent and obstacle, d_{\min} is the minimum safe distance, and \dot{d} is the derivative of distance with respect to time, $\xi = [\sigma, n, k]$ are the tunable parameters we desire to optimize in the online synthesis algorithm. The observations in Safety Gym include Lidar, speedometer, and magnetometer, which can be used to compute d and \dot{d} from observations. We compare the proposed algorithm against two types of baseline algorithms:

- CRL baselines, including TRPO-Lagrangian, PPO-Lagrangian and CPO [Achiam et al. \(2017\)](#); [Ray et al. \(2019\)](#). The cost threshold is set at zero to learn solid safe policies.
- FAC with original safety index ϕ_0 and handcrafted safety index ϕ_h , where $\phi_0 = d_{\min} - d$ and $\phi_h = 0.3 + d_{\min}^2 - d^2 - k\dot{d}$, named as *FAC with ϕ_0* and *FAC with ϕ_h* . The choice of ϕ_h is based on empirical knowledge. Details about baseline algorithms can be found in [Appendix C.1](#).

5.1. Evaluating FAC-SIS and Comparison Analysis

Results of the performance comparison are shown in [Figure 3](#). The results suggest that FAC with ϕ_0 performs poorest in the safety performance, which indicates that there are indeed many inevitably unsafe states, and SIS is necessary for these tasks. Only FAC-SIS learns the policy with *zero violation* and takes the lowest cost rate in all environments, answering the first question of zero constraint violation. FAC with ϕ_h fails to learn a solid safe policy in those environments with 0.30 constraint size (See the zoomed window in the second row), indicating that the handcrafted safety index can not cover all the environments. As for the baseline CRL algorithms, they can not learn a zero-violation policy in any environment because of the posterior penalty in the trial-and-error mechanism stated above. As for the reward performance, FAC-SIS has comparable reward performance in the `Push` task. For the `Goal` task, FAC with ϕ_h and FAC-SIS sacrifice the reward performance to guarantee safety, explained in [Ray et al. \(2019\)](#).

5.2. Validity Verification of Synthesized Safety Index

We conduct new metrics and experiments to show the validity of our synthesized safety index. Recall that the validity means that there always exists a feasible control policy to satisfy the constraint

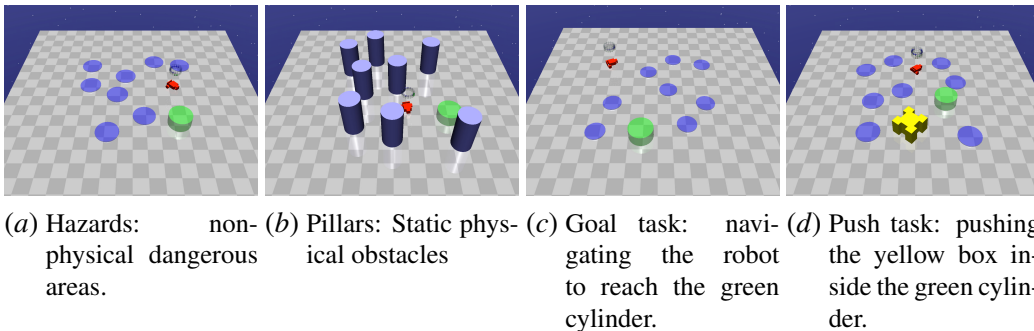


Figure 2: Obstacles and tasks in Safety Gym.

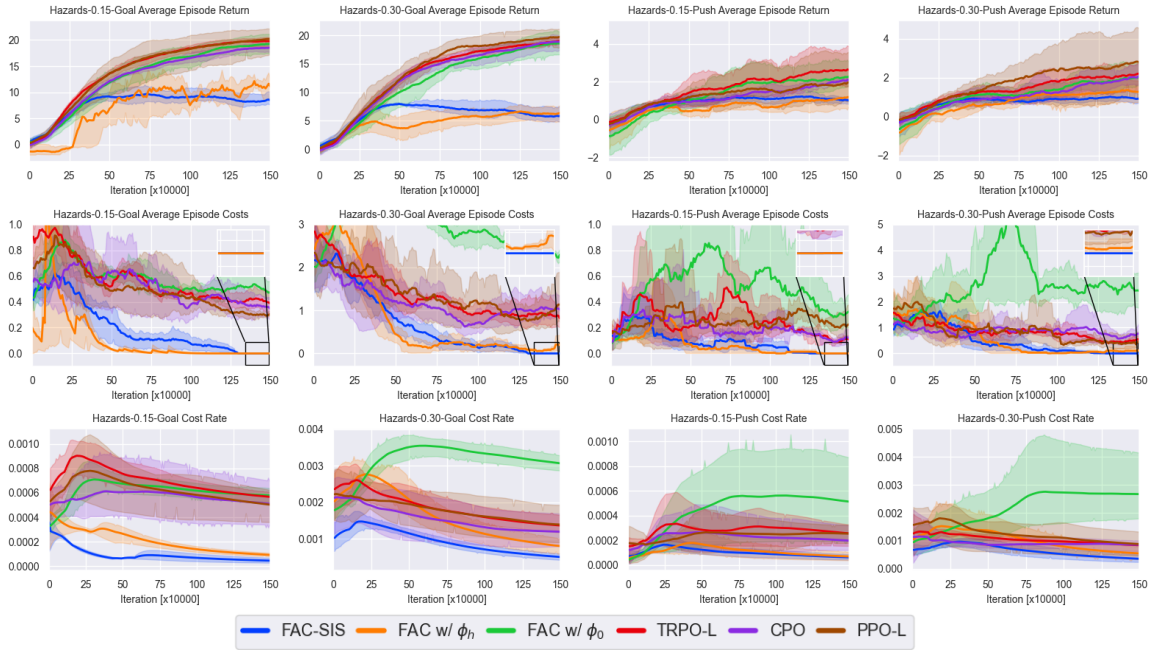


Figure 3: Average performance of FAC-SIS and baseline methods on 4 different Safety Gym environments over five seeds.

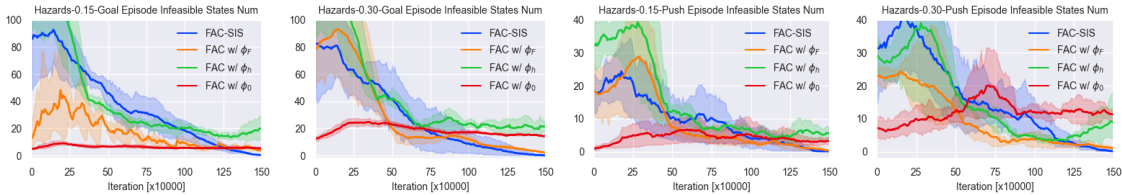


Figure 4: Average episodic number of violations of safe action constraint (1). A valid safety index and its corresponding safe control policy should have zero violation performance.

in (1). To effectively demonstrate the feasibility of SIS, we add a new baseline, *FAC with ϕ_F* , where ϕ_F is a valid safety index verified by Zhao et al. (2021). Figure 4 demonstrates the episodic number of constraint violations of (1) in the Safety Gym environments. The results show that FAC-SIS and *FAC with ϕ_F* can reach a nearly stable zero violation, which means that FAC can satisfy safe action constraint for a given valid safety index, and FAC-SIS also synthesizes a valid safety index. However, with ϕ_0 and ϕ_h there are consistent violations even with the converged policy, caused by different reasons. For ϕ_h , the reason is simply the inability to make the energy dissipate. For ϕ_0 , no high-order derivative in the safety index, so ϕ_0 cannot handle the high relative-degree between the constraint function and control input.

Furthermore, we want to give some scalable analyses about SIS. Firstly, we want to visualize that we can always find actions to dissipate the energy in the sampled distributions after SIS. We

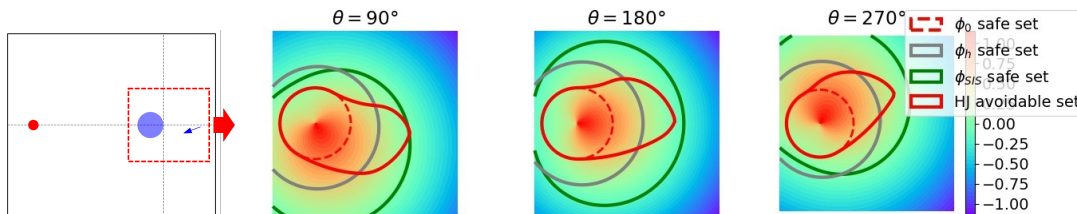


Figure 5: The color indicates the value of the synthesized safety index. The zero-sublevel set of learned safety index is a subset of the numerically max forward invariant set, the HJ avoidable set. We exclude all inevitably unsafe states through SIS like Figure 1. HJ-based solution is only capable for goal-reaching tasks with the reach-avoid problem formulation, while proposed algorithm can handle arbitrary task.

conduct experiments on a simple collision-avoidance environment shown in Figure 7³. We select three different initialization rules of the agent and hazard. We use a sample-based exhaustive method to identify infeasible states and see if they overlap the sampled state distributions. We project the state distribution to the 2D space of d and \dot{d} , and the results are listed in Figure 7. As we expect all the sampled states to be feasible, these two distributions *should not overlap*. The results show that the overlap of these two distributions is very small, which indicates that nearly all the sampling states are feasible.

Secondly, we want to visualize the shape of the safe set concerning the learned safety index. We slice the state space into three 2D planes with different agent heading angles, as shown in Figure 5. We use Hamilton-Jacobi reachability analysis to compute the avoidable sets numerically. The avoidable set considers the safe set under the *most conservative* control inputs, which is the maximum controlled invariant set (Mitchell, 2007; Choi et al., 2021). We use ϕ_h as the initial safety index of SIS. Figure 5 demonstrates that inevitable unsafe states exist in the zero-sublevel set of the empirical safety index. It also shows that we successfully exclude the inevitably unsafe sets through SIS. Notably, the zero-sublevel sets of the synthesized safety index are the subsets of the HJ-avoidable sets. The reasons why SIS can not learn the perfect shapes include limits of the representation capabilities of the safety index parameterization. Additionally, we still consider the *optimal criteria*, resulting in less conservative policies and possibly smaller safe sets.

6. Conclusion

This paper focuses on the joint synthesis of the safety certificate and the safe control policy for unknown dynamical systems and general tasks using a CRL approach. We are the first study to start with unknown dynamics and imperfect safety certificates, which significantly improves the applicability of the energy-function-based safe control. We add the optimization of safety index parameters as an outer loop of Lagrangian-based CRL with a unified loss. The convergence is analyzed theoretically. Experimental results demonstrate that the proposed FAC-SIS synthesizes a valid safe index while learning a safe control policy. In future work, we will consider more complex safety index parameterization rules, for example, neural networks. Meanwhile, we will consider other factors in SIS, such as the reward performance of the safe control policies.

3. See Appendix E.2 for more details.

Acknowledgments

This work was done during Haitong Ma’s internship at Carnegie Mellon University. This study is supported by National Key R&D Program of China with 2020YFB1600602 and Tsinghua University-Didi Joint Research Center for Future Mobility. This study is also supported by Tsinghua-Toyota Joint Research Fund. The authors would like to thank Mr. Weiye Zhao and Mr. Tianhao Wei for their valuable suggestions on the experiments.

References

- Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *International Conference on Machine Learning*, pages 22–31, Sydney, Australia, 2017. PMLR.
- Ayush Agrawal and Koushil Sreenath. Discrete control barrier functions for safety-critical control of discrete systems with application to bipedal robot navigation. In *Robotics: Science and Systems*, 2017.
- Aaron D Ames, Jessy W Grizzle, and Paulo Tabuada. Control barrier function based quadratic programs with application to adaptive cruise control. In *53rd IEEE Conference on Decision and Control*, pages 6271–6278. IEEE, 2014.
- Shalabh Bhatnagar and K Lakshmanan. An online actor–critic algorithm with function approximation for constrained markov decision processes. *Journal of Optimization Theory and Applications*, 153(3):688–708, 2012.
- Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural actor–critic algorithms. *Automatica*, 45(11):2471–2482, 2009.
- Vivek S Borkar. *Stochastic approximation: a dynamical systems viewpoint*, volume 48. Springer, 2009.
- Ya-Chien Chang, Nima Roohi, and Sicun Gao. Neural lyapunov control. *arXiv preprint arXiv:2005.00611*, 2020.
- Jianyu Chen, Shengbo Eben Li, and Masayoshi Tomizuka. Interpretable end-to-end urban autonomous driving with latent deep reinforcement learning. *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- Richard Cheng, Gábor Orosz, Richard M Murray, and Joel W Burdick. End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3387–3395, 2019.
- Jason J Choi, Donggun Lee, Koushil Sreenath, Claire J Tomlin, and Sylvia L Herbert. Robust control barrier-value functions for safety-critical control. *arXiv preprint arXiv:2104.02808*, 2021.
- Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained reinforcement learning with percentile risk criteria. *The Journal of Machine Learning Research*, 18(1):6070–6120, 2017.

- Jingliang Duan, Yang Guan, Shengbo Eben Li, Yangang Ren, Qi Sun, and Bo Cheng. Distributional soft actor-critic: Off-policy reinforcement learning for addressing value estimation errors. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596, Stockholm, Sweden, 2018. PMLR.
- Luis Gracia, Fabricio Garelli, and Antonio Sala. Reactive sliding-mode algorithm for collision avoidance in robotic systems. *IEEE Transactions on Control Systems Technology*, 21(6):2391–2399, 2013.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870, Stockholm, Sweden, 2018a. PMLR.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018b.
- Wanxin Jin, Zhaoran Wang, Zhuoran Yang, and Shaoshuai Mou. Neural certificates for safe control policies. *arXiv preprint arXiv:2006.08465*, 2020.
- Changliu Liu and Masayoshi Tomizuka. Control in a safe set: Addressing safety in human-robot interactions. In *Dynamic Systems and Control Conference*, volume 46209, page V003T42A003. American Society of Mechanical Engineers, 2014.
- Yuping Luo and Tengyu Ma. Learning barrier certificates: Towards safe reinforcement learning with zero training-time violations. *arXiv preprint arXiv:2108.01846*, 2021.
- Haitong Ma, Jianyu Chen, Shengbo Eben Li, Ziyu Lin, and Sifa Zheng. Model-based constrained reinforcement learning using generalized control barrier function. *arXiv preprint arXiv:2103.01556*, 2021a.
- Haitong Ma, Yang Guan, Shengbo Eben Li, Xiangteng Zhang, Sifa Zheng, and Jianyu Chen. Feasible actor-critic: Constrained reinforcement learning for ensuring statewise safety. *arXiv preprint arXiv:2105.10682*, 2021b.
- Paul Milgrom and Ilya Segal. Envelope theorems for arbitrary choice sets. *Econometrica*, 70(2): 583–601, 2002.
- Ian M Mitchell. A toolbox of level set methods. *UBC Department of Computer Science Technical Report TR-2007-11*, 2007.
- Stephen Prajna, Ali Jadbabaie, and George J Pappas. A framework for worst-case and stochastic safety verification using barrier certificates. *IEEE Transactions on Automatic Control*, 52(8): 1415–1428, 2007.
- Zengyi Qin, Kaiqing Zhang, Yuxiao Chen, Jingkai Chen, and Chuchu Fan. Learning safe multi-agent control with decentralized neural barrier certificates. *arXiv preprint arXiv:2101.05436*, 2021.

- Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement learning. *arXiv preprint arXiv:1910.01708*, 2019.
- Florian Richter, Ryan K Orosco, and Michael C Yip. Open-sourced reinforcement learning environments for surgical robotics. *arXiv preprint arXiv:1903.02090*, 2019.
- Ralph Tyrell Rockafellar. *Convex analysis*. Princeton university press, 2015.
- Ahmad EL Sallab, Mohammed Abdou, Etienne Perot, and Senthil Yogamani. Deep reinforcement learning framework for autonomous driving. *Electronic Imaging*, 2017(19):70–76, 2017.
- Matteo Saveriano and Dongheui Lee. Learning barrier functions for constrained motion planning with dynamical systems. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 112–119. IEEE, 2019.
- Mohit Srinivasan, Amogh Dabholkar, Samuel Coogan, and Patricio A Vela. Synthesis of control barrier functions using a supervised machine learning approach. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7139–7145. IEEE, 2020.
- Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning by pid lagrangian methods. In *International Conference on Machine Learning*, pages 9133–9143, Online, 2020. PMLR.
- Andrew Taylor, Andrew Singletary, Yisong Yue, and Aaron Ames. Learning for safety-critical control with control barrier functions. In *Learning for Dynamics and Control*, pages 708–717. PMLR, 2020.
- Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization. *arXiv preprint arXiv:1805.11074*, 2018.
- Eiji Uchibe and Kenji Doya. Constrained reinforcement learning from intrinsic and extrinsic rewards. In *2007 IEEE 6th International Conference on Development and Learning*, pages 163–168, Lugano, Switzerland, 2007. IEEE.
- Li Wang, Aaron D Ames, and Magnus Egerstedt. Safety barrier certificates for collisions-free multirobot systems. *IEEE Transactions on Robotics*, 33(3):661–674, 2017.
- Tianhao Wei and Changliu Liu. Safe control algorithms using energy functions: A unified framework, benchmark, and new directions. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 238–243. IEEE, 2019.
- Peter Wieland and Frank Allgöwer. Constructive safety using control barrier functions. *IFAC Proceedings Volumes*, 40(12):462–467, 2007.
- Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. *arXiv preprint arXiv:1905.11881*, 2019.
- Weiye Zhao, Tairan He, and Changliu Liu. Model-free safe control for zero-violation reinforcement learning. In *5th Annual Conference on Robot Learning*, 2021.

Appendix A. Theoretical Results in Section 3

A.1. Proof of Lemma 6

Define set of states with safe actions as $\mathcal{S}_f = \{s | \mathcal{U}_s(s) \neq \emptyset\}$. $\Delta = \mathbb{E}_s\{-v^{\pi^*}(s)\}$ since $-v^{\pi}(s)$ is irrelevant with ϕ . For $s \notin \mathcal{S}_f$ (i.e., $(\phi(s') - \max\{\phi(s) - \eta_D, 0\}) > 0$), we know that $\lambda^*(s) \rightarrow \infty$ from Lemma 5. Then $\lambda^*(s)$ is clipped to λ_{\max} . Therefore, the Lagrange function (6) can be reformulated to

$$\begin{aligned}
 & \mathcal{L}'(\pi^*, \lambda^*, \phi) \\
 &= \mathbb{E}_s \left\{ -v^{\pi}(s) + \lambda^*(s) (\phi(s') - \max\{\phi(s) - \eta_D, 0\}) \right\} \Big|_{\pi=\pi^*} \\
 &= \mathbb{E}_s \left\{ \lambda^*(s) (\phi(s') - \max\{\phi(s) - \eta_D, 0\}) \right\} \Big|_{\pi=\pi^*} + \Delta \\
 &= \mathbb{E}_{s \notin \mathcal{S}_f} \left\{ \lambda^*(s) (\phi(s') - \max\{\phi(s) - \eta_D, 0\}) \right\} \Big|_{\pi=\pi^*} + \Delta \\
 &= \lambda_{\max} \mathbb{E}_{s \notin \mathcal{S}_f} \left\{ (\phi(s') - \max\{\phi(s) - \eta_D, 0\}) \right\} \Big|_{\pi=\pi^*} + \Delta \\
 &= \lambda_{\max} J(\phi) + \Delta
 \end{aligned}$$

Appendix B. Theoretical Results in Section 4

B.1. Gradients Computation

The objective function of updating the policy and multipliers is the Lagrange function (3). Using the framework of maximum entropy RL, the objective function of policy update is:

$$J_{\pi}(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}} \left\{ \mathbb{E}_{a_t \sim \pi_{\theta}} \left\{ \alpha \log(\pi_{\theta}(a_t | s_t)) - Q_w(s_t, a_t) + \lambda_{\xi}(s_t) Q_{\phi}(s_t, a_t) \right\} \right\}$$

The policy gradient with the reparameterized policy $a_t = f_{\theta}(\epsilon_t; s_t)$ can be approximated by:

$$\begin{aligned}
 G_{\theta} = \hat{\nabla}_{\theta} J_{\pi}(\theta) &= \nabla_{\theta} \alpha \log(\pi_{\theta}(a_t | s_t)) + \left(\nabla_{a_t} \alpha \log(\pi_{\theta}(a_t | s_t)) \right. \\
 &\quad \left. - \nabla_{a_t} (Q_w(s_t, a_t) - \lambda_{\xi}(s_t) Q_{\phi}(s_t, a_t)) \right) \nabla_{\theta} f_{\theta}(\epsilon_t; s_t)
 \end{aligned}$$

where $\hat{\nabla}_{\theta} J_{\pi}(\theta)$ represents the stochastic gradient with respect to θ , and $Q_{\phi}(s_t, a_t) = (\phi(s_{t+1}) - \max\{\phi(s_t) - \eta_D, 0\})$. Neglecting those irrelevant parts, the objective function of updating the multiplier network parameters ξ is

$$J_{\lambda}(\xi) = \mathbb{E}_{s_t \sim \mathcal{D}} \left\{ \mathbb{E}_{a_t \sim \pi_{\theta}} \left\{ \lambda_{\xi}(s_t) (Q_{\phi}(s_t, a_t) - d) \right\} \right\}$$

The stochastic gradient is

$$G_{\xi} = \hat{\nabla} J_{\lambda}(\xi) = Q_{\phi}(s_t, a_t) \nabla_{\xi} \lambda_{\xi}(s_t) \tag{12}$$

The objective function of updating the safety index parameters ζ is already discussed, so the gradients for ζ is

$$G_{\zeta} = \lambda_{\xi}(s_t) \nabla_{\zeta} \Delta \phi^{\pi_{\theta}}(s_t) \tag{13}$$

where $\Delta \phi^{\pi_{\theta}}(s_t) = (\phi(s_{t+1}) - \max\{\phi(s_t) - \eta_D, 0\})|_{\pi_{\theta}}$, we use a different notation from $Q_{\phi}(s_t, a_t)$ since we focus on different variables.

B.2. Proof of Theorem 10

Recall the overview of the total convergence proof:

1. First we show that each update of the multi-time scale discrete stochastic approximation algorithm $(\theta_k, \xi_k, \zeta_k)$ converges almost surely, but at different speeds, to the *stationary point* $(\theta^*, \xi^*, \zeta^*)$ of the corresponding continuous-time system.
2. By Lyapunov analysis, we show that the continuous time system is locally asymptotically stable at $(\theta^*, \xi^*, \zeta^*)$.
3. We prove that (θ^*, ξ^*) is the locally optimal solution, or a local saddle point for the CRL problems with local optimal safety index parameters ζ^* .

First, we introduce the important lemma used for the convergence proof:

Lemma 12 (Soft policy evaluation, Haarnoja et al. (2018b)) *The Q -function update will converge to the soft Q -function as the iteration number goes to infinity.*

Remark 13 *In stochastic programming, the error in policy improvement caused by Q -function is a fixed bias rather than random variables, resulting in that it will not affect the convergence as long as the error is bounded. Therefore, we assume that Q -function is fully updated here for simplicity. Otherwise, the proof will be wordy.*

Lemma 14 (Convergence of clipped SGD, Zhang et al. (2019)) *For a stochastic gradient descent problem of a continuous differentiable and $L_0 - L_1$ smooth (which means $\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|$) loss function $f(x)$ and its stochastic gradient $\hat{\nabla} f(x)$, If these conditions are satisfied*

1. $f(x)$ is lower bounded;
2. There exists $\tau > 0$, such that $\|\nabla \hat{f}(x) - \nabla f(x)\| \leq \tau$ almost surely;

then the update of stochastic gradient descent converges almost surely with finite iteration complexity.

Then we continue to finish the multi-timescale convergence:

Remark 15 *In each timescale, we prove two facts, (1) the stochastic error is bounded, so the clipped gradient descent will converge. (2) the converged point is a stationary point. (Bhatnagar et al., 2009; Bhatnagar and Lakshmanan, 2012; Chow et al., 2017).*

Timescale 1 (Convergence of θ update).

Bounded error. As the random variable G_θ depends on the state-action pair sampled from replay buffer, so in the following derivation, it is denoted as $G_\theta(s, a)$. For the SGD using sampled (s_k, a_k) at k^{th} step, the stochastic gradient is denoted by $G_{\theta_k}(s_k, a_k)$. The error term with respect to θ is computed by

$$\delta\theta_k = G_{\theta_k}(s_k, a_k) - \mathbb{E}_{s \sim f^{\pi_\theta}} \{ \mathbb{E}_{a \sim \pi_\theta} G_{\theta_k}(s, a) \} \quad (14)$$

Therefore, the error term is bounded by

$$\begin{aligned}
 & \|\delta\theta_k\|^2 \\
 & \leq 2 \|f^{\pi_\theta}(s)\pi(a|s)\|_\infty^2 \left(\|G_{\theta_k}(s, a)\|_\infty^2 + |G_{\theta_k}(s_k, a_k)|^2 \right) \\
 & \leq 6 \|f^{\pi_\theta}(s)\pi(a|s)\|_\infty^2 \|G_{\theta_k}(s, a)\|_\infty^2
 \end{aligned} \tag{15}$$

where f^{π_θ} is the state distribution density under π_θ . As we assume the state and action are sampled from a closed set and the continuity of neural network, the upper bound is valid. According to Lemma 14 and invoking Theorem 2 in Chapter 2 of Borkar's book Borkar (2009), the optimization converges to a fixed point θ^* almost surely (for given ξ, ζ).

Stationary point θ^* . Then we show that the fixed point θ^* is a stationary point using Lyapunov analysis. The analysis of the fastest timescale, θ update, is rather easy, but it is helpful for the similar analyses in the next two timescales. According to Borkar (2009), we can regard the stochastic optimization of θ as a stochastic approximation of a dynamic system for given ξ, ζ :

$$\dot{\theta} = -\nabla_\theta \mathcal{L}'(\theta, \xi, \zeta) \tag{16}$$

Proposition 16 consider a Lyapunov function for dynamic system (16):

$$L_{\xi, \zeta}(\theta) = \mathcal{L}'(\theta, \xi, \zeta) - \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta) \tag{17}$$

where $\theta^*(\xi, \zeta)$ is a local minimum for given ξ, ζ . In order to show that θ^* is a stationary point, we need

$$dL_{\xi, \zeta}(\theta)/dt \leq 0 \tag{18}$$

Proof We have

$$\frac{dL_{\xi, \zeta}(\theta)}{dt} = -\|\nabla_\theta \mathcal{L}'(\theta, \xi, \zeta)\|^2 \leq 0 \tag{19}$$

The equality holds only when $\mathcal{L}'(\theta, \xi, \zeta) = 0$.⁴ ■

Combined with the conclusion of convergence, $\{\theta_k\}$ converges almost surely to a local minimum point θ^* for given ξ .

Timescale 2 (Convergence of λ update).

Bounded Error.

The error term of the ξ update,

$$G_{\xi_k}(s_k, a_k) - \mathbb{E}_s \{ \mathbb{E}_a G_{\xi_k}(s, a) \} \tag{20}$$

includes two parts:

4. Similar convergence proof in Chow et al. (2017) assumes that $\theta \in \Theta$ is a compact set, so they spend lots of effort to analyze the case when the θ reaches the boundary of Θ . However, the proposed clipped SGD has released the requirements of compact domain of θ , so the Lyapunov analysis becomes easier for the first timescale.

1. $\delta\theta_\epsilon$ caused by inaccurate update of θ (θ should converge to $\theta^*(\xi, \zeta)$ in Timescale 1, but to θ_k near $\theta^*(\xi, \zeta)$):

$$\begin{aligned}\delta\theta_\epsilon &= \hat{\nabla}_\xi \mathcal{L}'(\theta_k, \xi, \zeta) - \hat{\nabla}_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta) \\ &= (Q^{\pi_{\theta_k}}(s_k, a_k) - Q^{\pi_{\theta^*}}(s_k, a_k)) \nabla_\xi \lambda_\xi(s_k) \\ &= (\nabla_a Q(s_k, a_k) \nabla_\theta \pi(s_k) \epsilon_{\theta_k} + o(\|\epsilon_{\theta_k}\|)) \nabla_\xi \lambda_\xi(s_k)\end{aligned}\quad (21)$$

Therefore, $\|\delta\theta_\epsilon\| \rightarrow 0$ as $\|\epsilon_\theta\| \rightarrow 0$. The error is bounded since $\|\epsilon_\theta\|$ is a small error, where there must exist a positive scalar ϵ_0 s.t. $\|\epsilon_\theta\| \leq \epsilon_0$.

2. $\delta\xi_k$ caused by estimation error of ξ :

$$\delta\xi_k = G_{\xi_k}(s_k, a_k) - \mathbb{E}_{s \sim f^{\pi_\theta}} \{ \mathbb{E}_{a \sim \pi_\theta} G_{\xi_k}(s, a) \} \quad (22)$$

$$\|\delta\xi_k\|^2 \leq 4 \|f^{\pi_\theta}(s) \pi(a|s)\|_\infty^2 \left(\max \|Q_\phi(s, a)\|^2 + d^2 \right) \|\nabla_\xi \lambda_\xi(s_t)\|_\infty^2 \quad (23)$$

Similar to the analysis of Timescale 1 with compact domain of s_k , we can get the valid upper bound.

We again use Lemma 14 and Theorem 2 in Chapter 6 in [Borkar \(2009\)](#) to show that the sequence $\{\xi_k\}$ converges to the solution of the following ODE:

$$\dot{\xi} = -\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta) \quad (24)$$

Stationary point. Then we show that the fixed point ξ^* is a stationary point using Lyapunov analysis. Note that we have to take ϵ_θ into considerations.

Proposition 17 *For the dynamic system with the error term*

$$\dot{\xi} = -\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta) + \epsilon_\theta, \xi, \zeta) \quad (25)$$

Define a Lyapunov function to be

$$L_\zeta(\xi) = \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta) - \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) \quad (26)$$

where ξ^* is a local maximum point. Then $\frac{dL_\zeta(\xi)}{dt} \leq 0$.

Proof The proof is similar to Proposition 2; only the error of θ should be considered. We prove that the error of θ does not affect the decreasing property here:

$$\begin{aligned}\frac{dL_\zeta(\xi)}{dt} &= -(\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta) + \epsilon_\theta, \xi, \zeta))^T \nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta) \\ &= -\|\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta)\|^2 - \delta\theta_\epsilon^T \nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta) \\ &\leq -\|\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta)\|^2 + K_1 \|\epsilon_\theta\| \|\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta)\|\end{aligned}\quad (27)$$

where $K_1 = \|\nabla_a Q(s_k, a_k) \nabla_\theta \pi(s_k) \epsilon_{\theta_k}\|_\infty \|\nabla_\xi \lambda_\xi(s_k)\|_\infty \leq \infty$ according to (21) since the compact domain of s_k, a_k according to Assumption 2. As θ converges much faster than ξ according to the multiple timescale convergence in [Borkar \(2009\)](#), we get $dL_\zeta(\xi)/dt \leq 0$. Therefore, there exists trajectory $\xi(t)$ converges to ξ^* if initial state ξ_0 starts from a ball \mathcal{D}_{ξ^*} around ξ^* according to the asymptotically stable systems. \blacksquare

Local saddle point of (θ^*, ξ^*) . One side of the saddle point, $\mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) \leq \mathcal{L}'(\theta, \xi^*(\zeta), \zeta)$ are already provided in previous, so we need to prove here $\mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) \geq \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta)$. To complete the proof we need that

$$Q_\phi(s, a) \leq d \text{ and } \lambda^*(s)(Q_\phi(s, a) - d) = 0 \quad (28)$$

for all s in \mathcal{S}_f , and a sampled from $\pi_{\theta^*}(\xi, \zeta)$. Recall that λ^* is a local maximum point, we have

$$\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) = 0 \quad (29)$$

Assume there exists s_t and action a_t sampled from $\pi^*(s_t)$ so that $Q_\phi(s_t, a_t) > 0$. Then for λ^* we have

$$\nabla_\xi \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) = d_{\gamma^{\theta^*}}(s_t) Q_\phi(s_t, a_t) \nabla_\xi \lambda_\xi(s_t) \neq 0 \quad (30)$$

The second part only requires that $\lambda^*(s_t) = 0$ when $Q_\phi(s_t, a_t) < 0$. Similarly, we assume that there exists s_t and ξ^* where $\lambda_{\xi^*}(s_t) > 0$ and $Q_\phi(s_t, a_t) < d$. there must exists a ξ_0 subject to

$$\xi_0 = \xi^* + \eta_0 f^{\pi_{\theta^*}}(s_t) (Q_\phi(s_t, a_t) - d) \nabla_\xi \lambda_\xi(s_t) \quad (31)$$

for any $\eta \in (0, \eta_0]$ where $\eta_0 \geq 0$. It contradicts the statement the local maximum ξ^* . Then we get

$$\begin{aligned} & \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) \\ &= J_r(\theta^*) + \mathbb{E}_s \{ \lambda_{\xi^*}(s) \mathbb{E}_{a \sim \pi} \{ Q_\phi(s, a) \} \} = J_r(\theta^*) \\ &\geq J_r(\theta^*) + \mathbb{E}_s \{ \lambda_\xi(s) \mathbb{E}_{a \sim \pi} \{ Q_\phi(s, a) \} \} \\ &= \mathcal{L}'(\theta^*(\xi, \zeta), \xi, \zeta) \end{aligned} \quad (32)$$

So (θ^*, ξ^*) is a locally saddle point for given safety index parameters ζ .

Timescale 3 (Convergence of ζ update).

Bounded error. Similar to Timescale 2, the error of ζ update includes two parts

1. $\delta\theta_\epsilon + \delta\xi_\epsilon$ caused by inaccurate update of θ, ξ :

$$\begin{aligned} & \delta\theta_\epsilon + \delta\xi_\epsilon \\ &= \lambda_{\xi_k}(s_k) \nabla_\xi \Delta \phi^{\pi_{\theta_k}}(s_k) - \lambda_{\xi^*}(s_k) \nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s_k) \\ &= \lambda_{\xi_k}(s_k) \nabla_\zeta \Delta \phi^{\pi_{\theta_k}}(s_k) - \lambda_{\xi_k}(s_k) \nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s_k) \\ &\quad + \lambda_{\xi_k}(s_k) \nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s_k) - \lambda_{\xi^*}(s_k) \nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s_k) \\ &= (\lambda_{\xi_k}(s_k) - \lambda_{\xi^*}(s_k)) \nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s_k) \end{aligned} \quad (33)$$

The first part after the second equal sign is neglected since θ permutation does not affect on the gradient of ζ . Similar to derivation in (21), we get $|\delta\theta_\epsilon + \delta\xi_\epsilon| \rightarrow 0$ as $(\theta, \xi) \rightarrow (\theta^*, \xi^*)$.

2. $\delta\zeta_k$ caused by estimation error of ζ :

$$\delta\zeta_k = \lambda_{\xi^*}(s_k) \nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s_k) - \mathbb{E} \{ \lambda_{\xi^*}(s) \nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s) \} \quad (34)$$

The bounded error can be obtained by

$$\begin{aligned} & \|\delta\zeta_k\|^2 \\ &\leq 4 \|f^{\pi_{\theta^*}}(s) \pi(a|s)\|_\infty^2 \max_{s \in \mathcal{C}_D \mathcal{S}_f} |\lambda_{\xi^*}(s)|^2 \|\nabla_\zeta \Delta \phi^{\pi_{\theta^*}}(s)\|^2 \end{aligned} \quad (35)$$

Therefore, the ζ -update is a stochastic approximation of the continuous system $\zeta(t)$, described by the ODE For the dynamic system

$$\dot{\zeta} = -\nabla_{\zeta} \mathcal{L}'(\theta, \xi, \zeta)|_{\theta=\theta^*(\xi, \zeta)+\epsilon_{\theta}, \xi=\xi^*(\zeta)+\epsilon_{\xi}} \quad (36)$$

Stationary point. Define a Lyapunov function

$$L(\zeta) = \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) - \mathcal{L}'(\theta^*, \xi^*, \zeta^*) \quad (37)$$

where $\zeta^* \in Z$ is a local minimum point. Then

$$\begin{aligned} & \frac{dL(\zeta)}{dt} \\ &= (-\nabla_{\zeta} \mathcal{L}'(\theta^*(\xi, \zeta) + \epsilon_{\theta}, \xi^*(\zeta) + \epsilon_{\xi}, \zeta))^T \nabla_{\zeta} \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta) \\ &\leq -\|\nabla_{\zeta} \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta)\|^2 \\ &\quad + K_2 \|\epsilon_{\xi}\| \|\nabla_{\zeta} \mathcal{L}'(\theta^*(\xi, \zeta), \xi^*(\zeta), \zeta)\| \end{aligned} \quad (38)$$

where $K_2 = \|\nabla_{\xi} \lambda_{\xi}(s)\|_{\infty} \|\nabla_{\zeta} \Delta \phi^{\pi_{\theta^*}}(s)\|_{\infty}$ according to (33). The derivation is very similar to the conclusion in Proposition 17; the upper bound is valid since the compact domain of state and action. As ξ converges faster than ζ , $dL_{\zeta}(\xi)/dt \leq 0$, so there exists trajectory $\xi(t)$ converges to ξ^* if initial state ζ_0 starts from a ball \mathcal{D}_{ζ^*} around ζ^* according to the asymptotically stable systems.

Finally, we can conclude that the sequence $(\theta_k, \xi_k, \zeta_k)$, will converge to a locally optimal policy and multiplier tuple (θ^*, ξ^*) for a locally optimal safety index parameters, ζ^* .

Appendix C. Implementation Details

C.1. Codebase and Platforms

Implementation of FAC-SIS, FAC with ϕ_h and ϕ_F are based on the Parallel Asynchronous Buffer-Actor-Learner (PABAL) architecture proposed by Duan et al. (2021).⁵ All experiments are implemented on Intel Xeon Gold 6248 processors with 12 parallel actors, including 4 workers to sample, 4 buffers to store data and 4 learners to compute gradients. Implementation of other baseline algorithms are based on the code released by Ray et al. (2019)⁶ and also a modified version of PPO⁷.

C.2. Baseline Algorithms

The only difference between baseline algorithms, FAC with ϕ_h , ϕ_F is that no ζ update step in Algorithm 1.

Appendix D. Hyperparameters

The neural network design and detailed hyperparameters are listed in Table 4.

5. https://github.com/mahaitongdae/Safety_Index_Synthesis

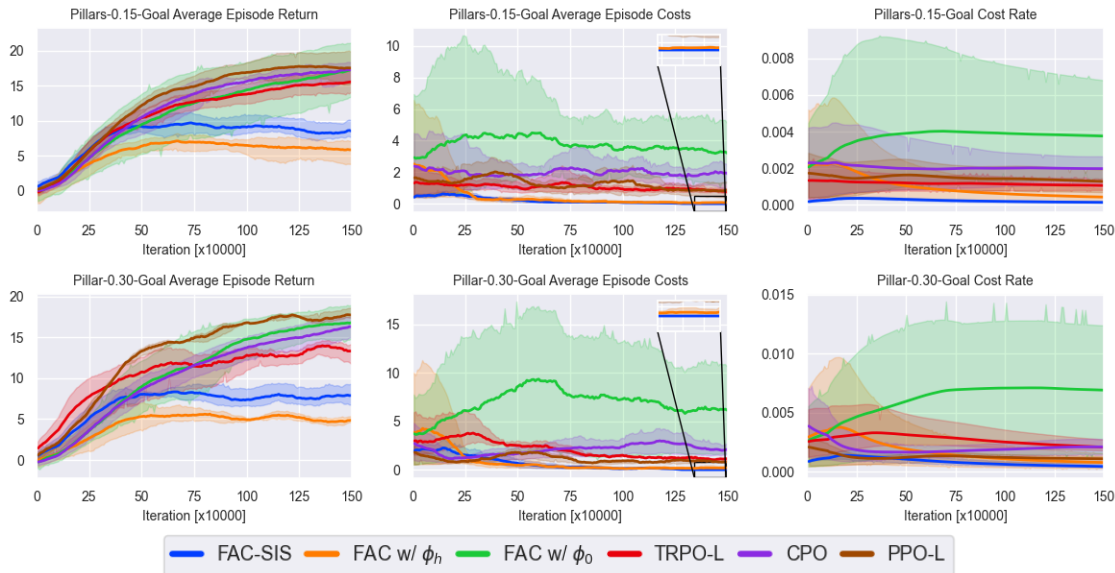
6. <https://github.com/openai/safety-starter-agents>

7. <https://github.com/ikostrikov/pytorch-a2c-ppo-acktr-gail>

Appendix E. Additional Experimental Results

E.1. Experiments in Other Safety Gym environments

We select six different Safety Gym environments, and the results of four of them are listed in the experiment section. The results with the rest of the environments are demonstrated here:



(a)

Figure 6: Average performance of FAC-SIS and baseline methods on 2 additional Safety Gym environments over five seeds.

The results in the other two environments are consistent with the experiment section. The proposed FAC-SIS learns a safe policy with zero constraint violation, and other baseline algorithms all fail to neglect the cost even in the converged policies. Furthermore, the reward performance is better than the handcrafted safety index, or FAC with ϕ_h .

E.2. Custom Environment Details

To effectively scale the state distribution, we manually set an environment similar to the Safety Gym environments with Point robot, Goal task and Hazard obstacles shown in Figure 7. The agent represented by the arrow should head to the red dot on the top while avoiding the hazard randomly located near the origin. The custom environment includes a static goal point at $(0, 5)$, a hazard with a radius of 0.5, and a point agent with two inputs, rotation and acceleration, the arrow represents the positive direction of the acceleration. The random initial positions (including positions, heading angles of the agent and the positions of hazard) design of three different distributions are listed in Table 1. The reward design includes two parts that are the tracking error of the heading angle towards goal position, and speed relevant to the distance to the goal position (It requires that the agent always takes 5 seconds to reach the goal, so the closer the agent get to the goal, the slower its target speed is.). The method to locate infeasible states is explained as follows. First, we discretize

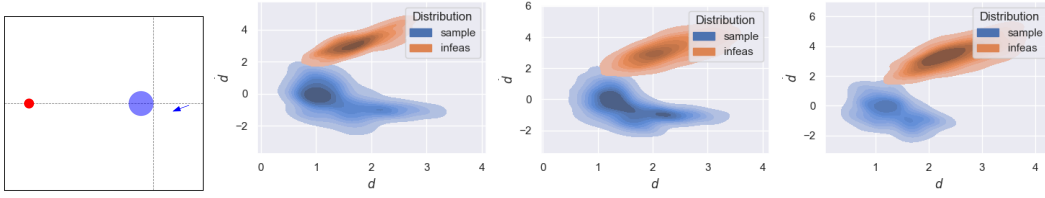


Figure 7: The custom environment and distributions of sampling and infeasible states under three different initialization setups. The overlap of the two distributions is very small, indicating that safe control exists for almost all sampled states.

the state and action space with small intervals, and we exhaust all the actions to see if the energy can dissipate for each state. If not, we remark the specific state as infeasible.

Distribution Index	Agent Initial Position			Hazard Initial Position	
	x	y	Angle	x	y
1	0	$[-1.5, -1.0]$	$[-\pi/4, \pi/4]$	0	$[0.5, 1]$
2	0	$[-1.5, -0.5]$	$[-\pi/4, \pi/4]$	0	$[0.5, 1, 5]$
3	$[-0.5, 0.5]$	$[-1.5, -0.5]$	$[-\pi/4, \pi/4]$	0	$[0.5, 1]$

Table 1: Initial distribution in custom environments. The initial position is sampled with a uniform distribution between the interval in the table.

E.3. Additional Results for Safety Index Synthesis

Parameters	Notion	k	σ	n
Synthesized	ϕ_{SIS}	0.7821	0.0958	1.149
Handcrafted	ϕ_h	1	0.3	2
Feasible	ϕ_F	1	0.04	2
Zero	ϕ_0	0	0	1

Table 2: Parameterization of Different Safety Index.

We list the different safety index parameters in TABLE 2. To be more specific, FAC-SIS has a single-direction update of the parameters by reducing the k, σ and changing n as shown in Figure 8. The trend can be explained by these two cases:

- $\phi(s) \leq 0$. Then the constraints, or the violation part are

$$\phi(s') = \sigma + d_{min}^n - d^n - kd \leq 0 \quad (39)$$

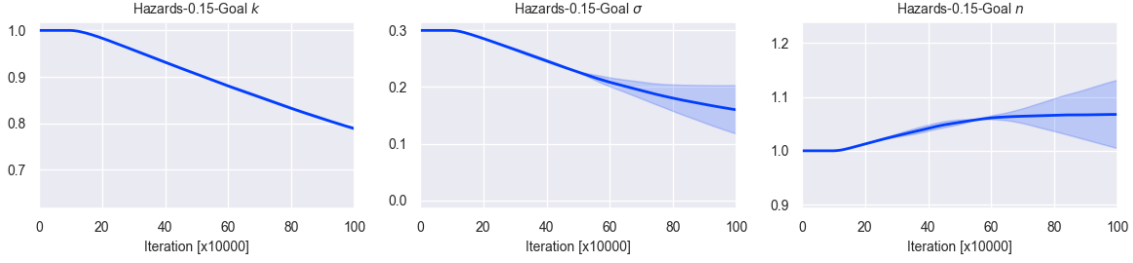


Figure 8: Learning curves of safety index parameters.

Therefore, if we want to further optimize the violation part J_ϕ , or the sum of LHS of the inequality in (39), then we have:

$$\begin{aligned} \partial J_\phi / \partial \sigma &= 1 > 0 \\ \partial J_\phi / \partial n &= n(d_{min}^{n-1} - d^{n-1}) < 0 \\ \partial J_\phi / \partial k &= -\dot{d} \end{aligned} \quad (40)$$

So, we should reduce k , increase n . The trend of k depends on the \dot{d} .

- $\phi(s) > 0$. If we further let $\eta_D = 0$, then we have $\phi(s') \leq \phi(s)$. The violation part of this inequality constraints are:

$$d^n - d'^n + kd\dot{d} - kd'\dot{d}' \quad (41)$$

we have

$$\begin{aligned} \partial J_\phi / \partial \sigma &= 0 \\ \partial J_\phi / \partial n &= n(d^{n-1} - d'^{n-1}) \\ \partial J_\phi / \partial k &= \dot{d} - \dot{d}' \end{aligned} \quad (42)$$

As the inequality constraints are violated, so at least there exists one positive term of $(d - d')$ and $(\dot{d} - \dot{d}')$. In other words, if we want to reduce the violation part (41), whether that one of the k or n reduces, or they all reduce.

Therefore, the synthesis trends in Figure 8 is reasonable.

Metric	ϕ_0	ϕ_i	ϕ_{SIS}	ϕ_F	ϕ_h
Success rate	0%	85%	99%	100%	70%
ϕ_0 violation rate	100%	0%	0%	0%	0%
Infeasible rate	100%	15%	1%	0%	30%
Average Tracking Error	-	2.974	3.378	3.456	3.528

Table 3: Performance Comparison of Different Safety Indexes.

We add quantified metrics on the custom environment to compare the feasibility, safety, and optimality of different safety indexes shown in TABLE 3. We randomly simulate 100 trajectories to see if the safety index leads to the infeasibility of unsafe actions. If the agent does not violate ϕ_0 (or stepping into hazards) or has infeasible states, then the trajectory is successful. ϕ_i is the initial safety index before synthesizing, and ϕ_{SIS} is the synthesized safety index. They both ensure safety but

the ϕ_{SIS} increases the feasibility by the safety index synthesis. The 1% difference might be caused by the mismatch of the state distributions between the custom environment and RL environments. ϕ_F indeed guarantees the best feasibility from the synthesis rules in [Zhao et al. \(2021\)](#). Besides, the synthesized safety index also has slightly better optimality, although we do not intend to improve it. This could be explained by the synthesized safety index being less conservative since it only ensures the sampled region by RL is feasible but not all the state space like the synthesis rules in [Zhao et al. \(2021\)](#).

Algorithm	Value
<i>FAC-SIS, FAC w/ ϕ_h, FAC w/ ϕ_F</i>	
Optimizer	Adam ($\beta_1 = 0.9, \beta_2 = 0.999$)
Approximation function	Multi-layer Perceptron
Number of hidden layers	2
Number of hidden units per layer	256
Nonlinearity of hidden layer	ELU
Nonlinearity of output layer	linear
Actor learning rate	Linear annealing $3e-5 \rightarrow 1e-6$
Critic learning rate	Linear annealing $8e-5 \rightarrow 1e-6$
Learning rate of multiplier net	Linear annealing $5e-6 \rightarrow 5e-6$
Learning rate of α	Linear annealing $8e-5 \rightarrow 8e-6$
Learning rate of safety index parameters (FAC-SIS only)	Linear annealing $8e-6 \rightarrow 1e-6$
Reward discount factor (γ)	0.99
Policy update interval (m_π)	3
Multiplier ascent interval (m_λ)	12
SIS interval (m_ϕ)	24
Target smoothing coefficient (τ)	0.005
Max episode length (N)	
Safety Gym task	1000
Custom task	120
Expected entropy ($\bar{\mathcal{H}}$)	$\bar{\mathcal{H}} = -\text{Action Dimensions}$
Replay buffer size	5×10^5
Replay batch size	256
Handcrafted safety index (ϕ_h) hyperparameters (η, n, k, σ)	(0, 2, 1, 0.3)
Feasible safety index (ϕ_F) hyperparameters (η, n, k, σ)	(0, 2, 1, 0.04)
<i>CPO, TRPO-Lagrangian</i>	
Max KL divergence	0.1
Damping coefficient	0.1
Backtrack coefficient	0.8
Backtrack iterations	10
Iteration for training values	80
Init λ	$0.268(\text{softplus}(0))$
GAE parameters	0.95
Batch size	2048
Max conjugate gradient iterations	10
<i>PPO-Lagrangian</i>	
Clip ratio	0.2
KL margin	1.2
Mini Bactch Size	64

Table 4: Detailed hyperparameters.