# Certificated Actor-Critic: Hierarchical Reinforcement Learning with Control Barrier Functions for Safe Navigation

Junjun Xie[1†], Shuhao Zhao[1†], Liang Hu[1*] and Huijun Gao[2*]

*Abstract*— **Control Barrier Functions (CBFs) have emerged as a prominent approach to designing safe navigation systems of robots. Despite their popularity, current CBF-based methods exhibit some limitations: optimization-based safe control techniques tend to be either myopic or computationally intensive, and they rely on simplified system models; conversely, the learning-based methods suffer from the lack of quantitative indication in terms of navigation performance and safety. In this paper, we present a new model-free reinforcement learning algorithm called Certificated Actor-Critic (CAC), which introduces a hierarchical reinforcement learning framework and well-defined reward functions derived from CBFs. We carry out theoretical analysis and proof of our algorithm, and propose several improvements in algorithm implementation. Our analysis is validated by two simulation experiments, showing the effectiveness of our proposed CAC algorithm.**

## I. INTRODUCTION

With the increasing deployment of autonomous robots and self-driving vehicles in the real world, the safety of autonomous navigation, e.g., collision avoidance with other users in shared spaces, has emerged as a common societal concern. Recently, the control barrier functions (CBFs) have been employed as a promising technique to address safety issues in robot navigation, either by being incorporated into traditional control and optimization frameworks or by being adapted into model-free learning approaches.

Building on control theory and optimization, many safe robot navigation approaches have been proposed. A typical paradigm is to divide navigation tasks into path planning and control [1], [2], with CBFs to ensure safety during path planning [3], [4], [5]. Another paradigm is to design the navigation controller directly [6], [7], in which classic control algorithms are combined with CBF in an optimization framework, such as quadratic programming (QP) [8], [9], [10] and model predictive control (MPC) [11]. However, QP-based methods are shortsighted due to pointwise optimization, resulting in locally optimal solutions or even a failed controller [12]. On the other hand, MPC mitigates the shortcomings of QP-based method but at the expense of a higher computational burden. In addition, the model-based safe control design usually requires a simplified explicit system model, which limits its application in navigating complex environments.

† Equal Contribution, * Corresponding Authors.

[1]J. Xie, S. Zhao and L. Hu are with the Department of Automation, School of Mechanical Engineering and Automation, Harbin Institute of Technology, Shenzhen, China.

[2]H. Gao is with Research Institute of Intelligent Control and Systems, Harbin Institute of Technology, Harbin, China.
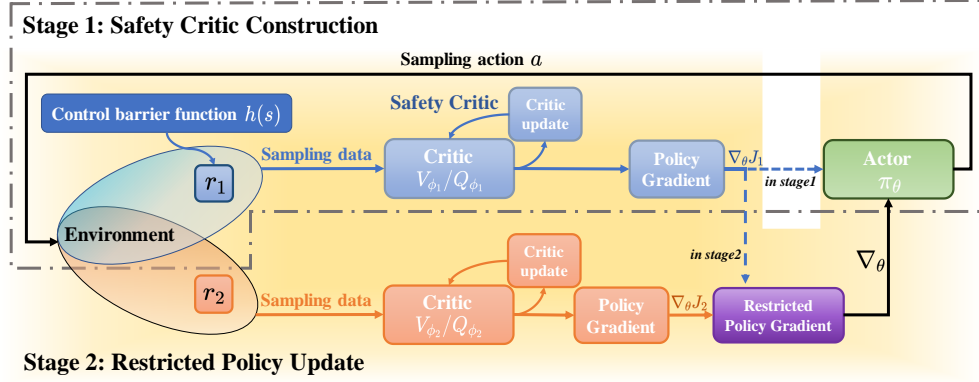
Reinforcement learning (RL) has long been a popular technique for robot navigation, particularly in scenarios where a model-free approach is required. Recently, researchers have explored various strategies to integrate CBFs with RL to enhance navigation safety [13], [14], [15]. Some methods impose CBFs as constraints on the learnt policy via a non-invasive approach [16], [17], [18]. That is, CBFs are used as a "safety filter" to adjust the trained policy via QP with CBF constraints. This approach is easy for implementation and has the merits of rigorous safety guarantees, but the adjustment on the origin policy might degrade its performance in an unpredictable approach. Another alternative direction is reward shaping using CBF [19], [20] in which safety is combined together with other reward terms. Since the final reward function is a trade-off between safety and other objectives, the safety requirement cannot be guaranteed, and it lacks quantitative analysis and explainability in terms of safety.

To overcome the limitation in existing methods, we propose a novel model-free reinforcement learning algorithm called **Certificated Actor-Critic**. Our algorithm features a CBFs-derived reward function that provides a quantitative assessment of navigation safety. Furthermore, we develop a hierarchical RL framework that learns a safe policy first and then refines it at the next stage to achieve fast goal-reaching without compromising safety.

The main contributions of the paper are threefold:

1) We propose a model-free reinforcement learning algorithm certificated Actor-Critic including a CBFs-derived reward function, which can quantitatively estimate the safety of the policies and states;

2) We design a hierarchical framework that accommodates safety and goal-reaching objectives in robot navigation, and improve its goal-reaching capability yet maintaining safety via novel restricted policy update strategies;

3) We conduct two experiments with detailed comparative analysis, showing the effectiveness of our proposed algorithm.

## II. PRELIMINARIES

Consider an infinite-horizon Markov decision process (MDP) [21] with a discrete-time stochastic system

$$s_{t+1} = \mathcal{F}(s_t, a_t) \tag{1}$$

which can be defined concisely as a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{F}, R, \gamma \rangle$, where $\mathcal{S}$ is a set of states satisfied $s_t \in \mathcal{S} \subseteq \mathbb{R}^n$, $\mathcal{A}$ is a set of actions satisfied $a_t \in \mathcal{A} \subseteq \mathbb{R}^m$, $\mathcal{F} : \mathcal{S} \times \mathcal{A} \to \mathcal{S}$ is

Fig. 1. The framework of certificated actor-critic.

a possibly dynamics, $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is a reward function, and $\gamma \in [0, 1]$ is a discount factor.

### A. Reinforcement Learning

Reinforcement learning aims to learn a policy via interaction between environment and agent with the dynamic model (1). The cumulative reward is defined as $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k}$. With a policy $\pi$, state value function $V^\pi(s)$ and action-value function $Q^\pi(s, a)$ are defined as

$$
\begin{aligned}
V^\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s], \\
Q^\pi(s, a) &= \mathbb{E}_\pi[G_t | S_t = s, A_t = a].
\end{aligned}
\tag{2}
$$

The objective is to find an optimal policy $\pi^*$ to make $V^{\pi^*}(s) \geq V^\pi(s), \forall \pi, s \in \mathcal{S}$.

### B. Discrete-time Control Barrier Functions

Consider the discrete-time system (1), we have the following definitions and lemma:

**Definition 1.** *(Safe Set [22], [23], [24]) The **safe set** $\mathcal{C}$ is defined as the zero-superlevel set of a smooth function $h : \mathbb{R}^n \to \mathbb{R}$, i.e.,*

$$
\begin{aligned}
\mathcal{C} &= \{s \in \mathbb{R}^n : h(s) \geq 0\}, \\
\partial\mathcal{C} &= \{s \in \mathbb{R}^n : h(s) = 0\}, \\
\mathrm{Int}(\mathcal{C}) &= \{s \in \mathbb{R}^n : h(s) > 0\}.
\end{aligned}
\tag{3}
$$

**Definition 2.** *(Discrete-time Control Barrier Function [25]) The function $h : \mathbb{R}^n \to \mathbb{R}$ is called a **Discrete-time Control Barrier Function** for the system (1) if there exists $\alpha$ such that*

$$
\sup_{a_t \in \mathcal{A}} [h(\mathcal{F}(s_t, a_t)) + (\alpha - 1)h(s_t)] \geq 0, 0 < \alpha < 1.
\tag{4}
$$

*Lemma* 1. (Forward Invariant) The system (1) is **forward invariant** in safe set $\mathcal{C}$, i.e.,

$$
\forall s_t \in \mathcal{C} \Rightarrow s_{t+1} \in \mathcal{C}
\tag{5}
$$

if

$$
h(s_{t+1}) + (\alpha - 1)h(s_t) \geq 0, 0 < \alpha < 1.
\tag{6}
$$

*Proof.* Since $s_t \in \mathcal{C}$, from **Definition** 1 we have $h(s_t) \geq 0$. If (6) holds on, then $h(s_{t+1}) \geq (1 - \alpha)h(s_t) \geq 0$, which means $s_{t+1} \in \mathcal{C}$. ∎

---

**Algorithm 1** Certificated Actor-Critic

1: Design the control barrier function $h(s)$ for the system (1) with expected decay rate $\alpha_0$
2: *%Stage 1: Safety Critic Construction*
3: Set reward $r_1$ as (12), initialize the actor network $\pi_\theta$ with parameters $\theta$ and the critic network $Q_{\phi_1}$ or $V_{\phi_1}$ with parameters $\phi_1$
4: Define learning rate $\lambda_\theta, \lambda_{\phi_1}$, the loss function $J_1(\theta)$ for actor and $L_1(\phi_1)$ for critic with $r_1$
5: **for** each step in training **do**
6:     $\theta \leftarrow \theta - \lambda_\theta \nabla_\theta J_1(\theta)$
7:     $\phi_1 \leftarrow \phi_1 - \lambda_{\phi_1} \nabla_{\phi_1} L_1(\phi_1)$
8: **end for**
9: *%Stage 2: Restricted Policy Update*
10: Set reward $r_2$, initialize the critic network $Q_{\phi_2}$ or $V_{\phi_2}$ with parameters $\phi_2$ for navigation
11: Define learning rate $\lambda_{\phi_2}$, the loss function $J_2(\theta)$ for actor and $L_2(\phi_2)$ for critic with $r_2$
12: **for** each step in training **do**
13:     $\nabla_\theta \leftarrow$ (10)
14:     $\theta \leftarrow \theta - \lambda_\theta \nabla_\theta$
15:     $\phi_1 \leftarrow \phi_1 - \lambda_{\phi_1} \nabla_{\phi_1} L_1(\phi_1)$
16:     $\phi_2 \leftarrow \phi_2 - \lambda_{\phi_2} \nabla_{\phi_2} L_2(\phi_2)$
17: **end for**
18: **Output:** $\theta$, $\phi_1$ and $\phi_2$

---

## III. PROBLEM FORMULATION

In this section, we consider a robot navigation task, and formulate the problem as follows:

Given the robot dynamics (1) with a proper control barrier function $h(s)$, we aim to learn a navigation policy such that:

1) (*Safety Guaranteed*) The robot should stay in the safe set defined in **Definition** 1 with the corresponding control barrier function $h(s)$;
2) (*Goal reaching*) The robot should reach its navigational goal while keeping in its safe conditions;
3) (*Safety Certificates*) A quantitative safety certificate should be provided, illustrating the level of safety assured by the navigation policy via numerical values.

## IV. Certificated Actor-Critic

In this section, we propose a hierarchical reinforcement learning framework called ***Certificated Actor-Critic (CAC)*** that extends from the classic actor-critic architecture [21]. We decompose the entire navigation task into two sub-tasks that consider safety and goal-reaching respectively. Accordingly, two separate critics are used to update the actor in two stages, as shown in **Algorithm** 1 and Fig. 1.

Different from the common approach that combines multi-objectives into a single objective by weight coefficients, CAC considers safety and goal-reaching in two consecutive stages. In the first stage of **safety critic construction**, an initial policy is learnt with the reward of safety. Then in the second stage of **restricted policy update**, the obtained policy from the previous stage is further updated subject to both rewards of safety and that of goal-reaching, which improves goal-reaching performance yet without compromising safety. More significantly, due to the well-designed reward function, the critic trained based on rewards of safety also works as a safety certificate to evaluate safety of the learnt policy.

### A. Safety Critic Construction

In the first stage, we construct a safety critic and train a safe policy via a well-designed reward of safety.

Consider a system with a defined control barrier function $h(s)$ and expected decay rate $\alpha_0$. From *Lemma* 1, if the action $a_t$ satisfies (6) at a particular safe state $s_t$, then the state at the next step $s_{t+1} = f(s_t, a_t)$ is still safe. Otherwise, if (6) is violated, we have $h(s_{t+1}) + (\alpha_0 - 1)h(s_t) < 0$. Hence, we define

$$r_1(s_t, a_t) = \delta_h \triangleq \min(h(s_{t+1}) + (\alpha_0 - 1)h(s_t), 0) \quad (7)$$

and regard $\delta_h$ as a **safety certificate** for action $a_t$, i.e., if $\delta_h(s_t, a_t) = 0$, the system is safe at step $t$.

If every step in an episode satisfies $\delta_h = 0$, the system is safe during the whole process. Motivated by this, we call the state value function $V_1^\pi$ and the action value function $Q_1^\pi$ **safety critics** due to their ability on safety evaluation, formally stated in *Theorem* 1.

*Theorem* 1. **(Safety Critic)** Consider the system (1) with reward function (7):

1) The system (1) is safe with policy $\pi$ from the initial state $s_0$ if $s_0 \in \mathcal{C}$ and $V_1^\pi(s_0) = 0$;
2) The system (1) is safe with policy $\pi$ with the initial state-action pair $(s_0, a_0)$ if $s_0 \in \mathcal{C}$ and $Q_1^\pi(s_0, a_0) = 0$.

*Proof.* For 1), if $V_1^\pi(s_0) = 0$, we have that the return of any episode from $s_0$ satisfies

$$\sum_{k=0}^{\infty} \gamma^k r_1(s_k, a_k) = 0. \quad (8)$$

From (7), $r_1 = \min(h(s_{t+1}) + (\alpha_0 - 1)h(s_t), 0) \leq 0$, so the unique solution is $r_1(s_k, a_k) = 0, \forall k \geq 0$. From *Theorem* 1, $\forall k, s_k \in \mathcal{C}$, i.e., the system is always safe from $s_0$.

For 2), if $Q_1^\pi(s_0, a_0) = 0$, then the return of any episode choosing $a_0$ at initial state $s_0$ also satisfies (8). Similarly,
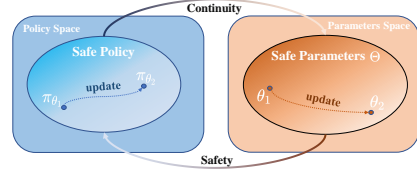


Fig. 2. The relationship between safe policy and its parameters under *Assumption* 1. Since the corresponding parameters of similar safe policies are continuous, parameters can be updated in a small neighbourhood, which guarantees the safety of policy.

since $r_1 \leq 0$, we have $\forall k, s_k \in \mathcal{C}$ from *Theorem* 1, and the system is always safe from $s_0$. ∎

According to *Theorem* 1, we can use safety critics: 1) to evaluate relative safety between policies; for example, if $V_1^{\pi_1}(s) \geq V_1^{\pi_2}(s)$ at all states, then $\pi_1$ is safer than $\pi_2$; 2) to determine the absolute safety about some states; if $v_\pi(s_0) = 0$, then the system is absolutely safe from the initial state $s_0$. It is also worth emphasizing that even the optimal safe policy $\pi^*$ does not necessarily satisfy $V_1^{\pi^*} = 0$ because of limited action space $\mathcal{A}$. For example, the robot is too close to avoid the obstacle, for which no irretrievable actions exist.

After training, we derive an optimal safe policy $\pi_{\text{safe}}^*$, safety critics $V_1^{\pi^*}$ and $Q_1^{\pi^*}$. What remains is to update the policy for better goal-reaching behaviour, which is left for the next stage of restricted policy update.

### B. Restricted Policy Update

The core challenge in this stage is how to guarantee the safety of the system when policy is updated for better goal-reaching performance. That is, the values of safety critic do not decrease for all states in the second stage. Here, we make an assumption based on empirical observations that safe states and actions are usually non-unique. For example, in the lane-keeping task, the car could drive safely not only in the middle of the lane but also on the left or right slightly. The assumption is stated as below:

*Assumption* 1. **(Parameter Continuity of Safe Policies)** Given a parameterized safe policy $\pi_\theta$, $\forall \delta > 0$, $\exists \theta'$ subject to $\|\theta' - \theta\| \leq \delta$, then the policy $\pi_{\theta'}$ is safe as well.

*Assumption* 1 implies that safe policies are not isolated but continuous in parametric space, and Fig. 2 shows the relationship between safe policy and safe parameters. It leaves us spaces to update the policy for better goal-reaching performance while keeping safe.

In the second stage, we can set any appropriate reward function $r_2(s_t, a_t)$ for goal-reaching in navigation. For example, set

$$r_2(s_t, a_t) = -\max(l(s_{t+1}) + (\beta_0 - 1)l(s_t), 0) \quad (9)$$

where $l(s) : \mathbb{R}^n \to \mathbb{R}$ is a control Lyapunov function (CLF) [26], [27] with an expected decay rate $\beta_0$. In the actor-critic framework, the critic provides the gradient $\nabla_\theta J(\theta)$ for actor (policy) to update. To keep safety critic do not decrease, we should determine the gradient depending on both safety critic

and navigation critic. Specifically, we derive the gradient $\nabla_\theta$ using **restricted policy update** as

$$\nabla_\theta = \arg\max_e \, e \cdot \nabla_\theta J_2(\theta)$$
$$\text{s.t.} \quad e \cdot \nabla_\theta J_1(\theta) \geq 0 \tag{10}$$
$$\|e\| \leq \|\nabla_\theta J_2(\theta)\|$$

where $J_1(\theta), J_2(\theta)$ are actor loss functions with reward $r_1, r_2$ respectively, and $\nabla$ is gradient operator. Under *Assumption* 1, the policy $\pi^*_{\text{safe}}$ will update and converge gradually to the final optimal policy $\pi^*$ along $\nabla_\theta$.

### C. Practical Improvements

So far, we have introduced the overall framework of our algorithm, which can be transplanted into any existing actor-critic architecture. However, for training in deep reinforcement learning where all actors and critics are represented as neural networks, we need to make a few improvements.

*1) Policy improvement:* Consider the parameterized actor network $\pi_\theta$, a widely used policy gradient method is to maximize the function $J(\theta) = \mathbb{E}[V^{\pi_\theta}(s)]$. Since the target is an expectation of state values, there may be cases that some values rise and others decline. On safety, it means some states become safer but others more dangerous. Hence, we improve the policy in the way soft actor-critic [28], [29] does, that is using Kullback-Leibler divergence as $J(\theta)$ leads to

$$\pi_{\text{new}} = \arg\min_{\pi' \in \Pi} D_{\text{KL}} \left( \pi'(\cdot \mid \mathbf{s}_t) \, \Big\| \, \frac{\exp\left(Q^{\pi_{\text{old}}}(\mathbf{s}_t, \cdot)\right)}{Z^{\pi_{\text{old}}}(\mathbf{s}_t)} \right) \tag{11}$$

and a significant conclusion is $Q^{\pi_{\text{new}}}(s,a) \geq Q^{\pi_{\text{old}}}(s,a)$ [28], [29]. It shows that all states become safer after policy update.

*2) Exponential reward normalization:* We define $r_1$ in (7), and the best reward at each step is $0$. However, in finite episode, it leads to early termination, since the longer the episode is, the more negative the cumulative reward is. So we adjust reward function to

$$r_1(s_t, a_t) = \exp(\min(h(s_{t+1}) + (\alpha_0 - 1)h(s_t), 0)) \tag{12}$$

and obviously $r_1 \in (0, 1]$. Besides, if there exists maximum length $k_{\max}$ of an episode, then $V_1^\pi \leq \frac{1-\gamma^{k_{\max}}}{1-\gamma}$, which means the critic network still can work as the safety critic with a bias $\frac{1-\gamma^{k_{\max}}}{1-\gamma}$, i.e., if $V_1^\pi(s_0) = \frac{1-\gamma^{k_{\max}}}{1-\gamma}$, then the system is safe from the initial state $s_0$.

*3) Gradient enhancement:* A restricted gradient is derived in (10) to guarantee safety performance during the second stage. Although $\nabla_\theta \cdot \nabla_\theta J_1(\theta) \geq 0$ works theoretically, there exist two reasons for possible failure in algorithm implementation. One is there always needs a step length to update the parameters, and local gradient does not guarantee global convergence. Another is the true gradient $\nabla_\theta J_1(\theta)$ is unknown, and is replaced by the estimation $\widetilde{\nabla}_\theta J_1(\theta)$ derived from data. Hence, it is essential to enhance the constraint as $\nabla_\theta \cdot \widetilde{\nabla}_\theta J_1(\theta) \geq \delta$ or $\cos\{\nabla_\theta, \widetilde{\nabla}_\theta J_1(\theta)\} \geq \delta, \delta > 0$ where $\cos\{\cdot, \cdot\}$ represents the cosine of angle between two vector.
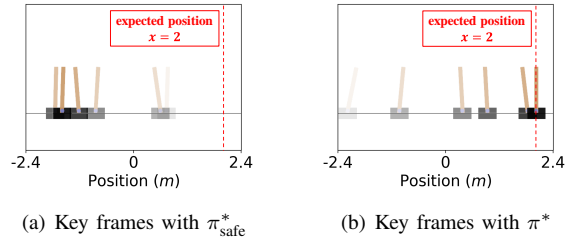


(a) Key frames with $\pi^*_{\text{safe}}$      (b) Key frames with $\pi^*$

Fig. 3. Frames in an episode with $\pi^*_{\text{safe}}$ and $\pi^*$.



(a) Positions with $\pi^*_{\text{safe}}$      (b) Angles with $\pi^*_{\text{safe}}$

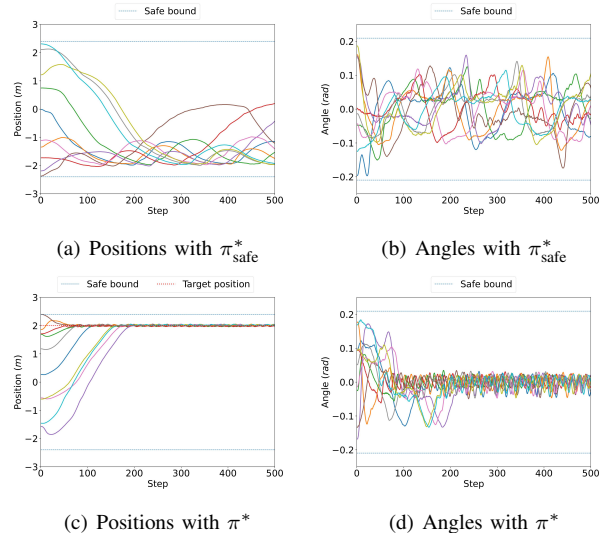(c) Positions with $\pi^*$      (d) Angles with $\pi^*$

Fig. 4. Positions and angles of CartPole with $\pi^*_{\text{safe}}$ and $\pi^*$ in 10 test episodes. Both policies guarantee the state in the safe range, and the final policy $\pi^*$ drives the CartPole to the target position in safe conditions.

## V. EXPERIMENTS

In this section, we validate our CAC algorithm on two simulation experiments: a classic control task CartPole using Gymnasium [30], and an autonomous underwater vehicle (AUV) navigation task using the simulator HoloOcean [31].

### A. Continuous CartPole

*1) Basic Experiment:* CartPole is a typical benchmark for reinforcement learning and control. The task is to balance the pole as long as possible by applying the force on the cart. The state is defined as $s = \begin{bmatrix} x & v & \theta & \omega \end{bmatrix}^T$, where $x, v$ are the position, velocity of the cart, and $\theta, \omega$ are the angle, angular velocity of the pole respectively. We set the environment with the continuous action space $a \in [-1, 1]$ and the default state space. The navigation requirements are

$$\text{Safety (allowable states): } -2.4 \leq x \leq 2.4$$
$$-12° \leq \theta \leq 12°$$

Navigation destination (desired states): $x = 2$

The initial state is sampled randomly from allowable states.

For safety guaranteed, we design two control barrier functions $h_1(s) = (2.4^2 - x^2)/2.4^2$ and $h_2(s) = (12^2 - \theta^2)/12^2$. From (7) and (12), the reward of the first stage is $r_1 = [\exp(\delta_{h_1}) + \exp(\delta_{h_2})]/2$.
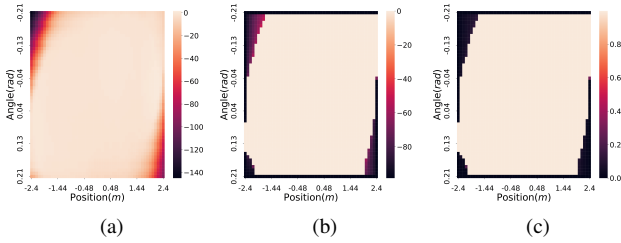
Fig. 5. Heatmaps of (a) safety critic value, (b) average sampling return and (c) safe rate. Three heatmaps are similar and consistent, which validates that the safety critic is a good safety certificate.
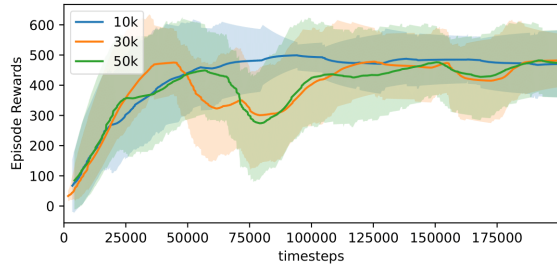


Fig. 6. Episode rewards on safety (maximum is 500) of the CAC models that are trained with different lengths at the first stage. All models perform almost the same in safety after 200k steps of training at stage 2.

*2) Verification of Safety Critic V:* As mentioned in *Theorem* 1, thanks to our well-defined reward, the critic network of the first stage represents the safety of the policy. To verify the conclusion, we compare the safety critic value $V_1(s)$ with the average sampling return and safe episodes rate (an episode is safe if its length is $k_{\max} = 500$) from 10 episodes with all states using $\pi^*_{\text{safe}}$. The results are depicted in Fig. 5. Obviously, the three heatmaps are consistent, which verifies the conclusion in *Theorem* 1.

Furthermore, we also study the effect of the first stage performance on the final results. Specifically, we respectively train 10k, 30k and 50k steps in the first stage, and 200k steps in the second. The policy networks with 10k and 30k steps do not fully converge in the first stage, but all three

final models after the second stage work well. It shows that even if the first stage is not fully trained, the CAC algorithm is still effective with sufficient training steps in the second stage. Moreover, we test the models in the process, which are trained 100k steps in the second stage, and there are significant improvements or maintenance of a high level in safety for all models. shown as Fig. 6 and Table I .

TABLE I
SAFE RATE IN 100 EPISODES

| Length of stage 1 | Stage 1 | After 100k steps of stage 2 | Stage 2 |
|---|---|---|---|
| 10k steps | 63% | 99% | 100% |
| 30k steps | 95% | 100% | 100% |
| 50k steps | 100% | 100% | 100% |



(a) Average navigation step reward    (b) Length of episodes

Fig. 8. Comparison of average navigation reward per step and length of training episodes in the underwater environment for the CAC models trained with and without policy restriction (10).

*B. Autonomous Underwater Vehicle*

To examine the safe navigation capability in complex environments with dense obstacles, we consider an Autonomous Underwater Vehicle (AUV) navigation task using the HoloOcean simulator [31].

In simulation, the HoveringAUV is kept at a fixed depth. The forward velocity is set to 0.9m/s and the yaw angular velocity range is (-1,1) rad/s. The velocities are followed via a simple PD controller. The AUV is equipped with a RangeFinder sensor with 120° FOV, whose measurement is a $\mathbb{R}^9$ vector $[d_1, \cdots, d_9]$ representing the distances with 15° interval (the maximum distance is 10 meters). The other
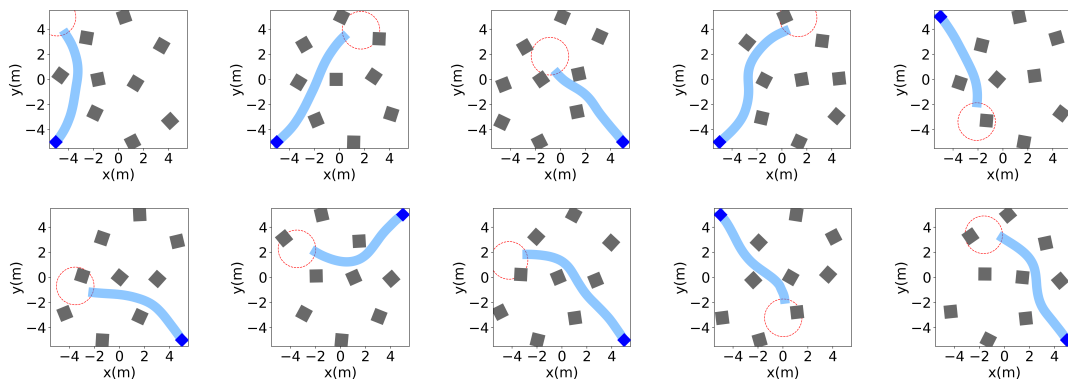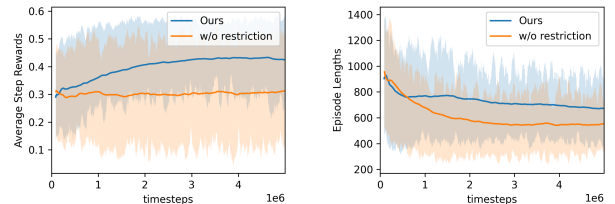


Fig. 7. Results of test trajectories (light blue). There are 9 random cubic obstacles (grey) in each case. The initial position (dark blue) is set as one of four corners of the pool, and target position (red circle) is generated in the diagonal half of the initial position randomly.
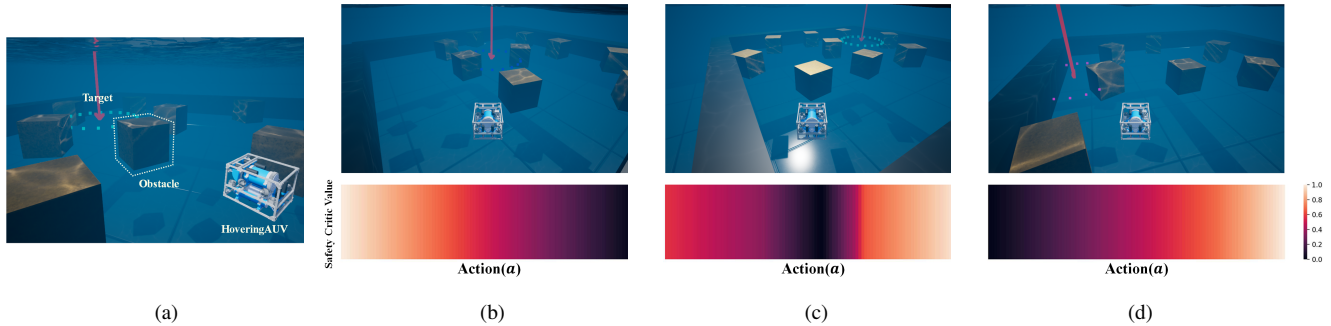
Fig. 9. (a): The simulation environment in HoloOcean. (b-d): Results of safety critic Q-values under particular observations. The horizontal coordinate of the heat map represents the action from -1 rad/s to 1 rad/s. The color shade represents the safety critic Q-value for performing the corresponding action under the current observation. The values are normalized to be distinguished obviously.

TABLE II

ABLATIONS RESULTS IN 100 RANDOM EPISODES

| Algorithm | Success Rate$^a$ ↑ | Collision Rate ↓ | Average length$^b$↓ |
|---|---|---|---|
| T-O(0.5)$^c$ | 48% | 31% | 852 |
| T-O(0.25) | 67% | 33% | 679 |
| Stage 1 | N/A | 18% | N/A |
| W/o Re.$^d$ | 49% | 51% | **616** |
| Our CAC | **86%** | **11%** | 625 |

a): reaching target position without collision is regarded as a successful episode;

b): the average length of successful episodes;

c): T-O($x$) represents the reward function is set as a trade-off with $x$ as the coefficient of safety reward, and $1-x$ as that of navigation reward;

d): policy update only using $\nabla_\theta J_2(\theta)$ without restriction (10).

observations are the yaw angular velocity $\omega_\gamma$ from the IMU sensor, horizontal velocities $[v_x, v_y]$, the location of AUV $[x, y]$, the rotation yaw angle $\gamma$ of AUV in the global frame and the goal location $[x_{tar}, y_{tar}]$. All observation vectors are concatenated into a $\mathbb{R}^{17}$ vector after normalization.

The simulation environment is a $10m \times 10m$ pool, as depicted in Fig. 9(a). The AUV is initialized randomly at one of the four corners of the pool. The obstacles, which are cubic, are randomly generated in the pool and distributed densely. The target position is randomly generated in the other diagonal half area of the pool where the starting point is not located. When the distance between the AUV and the target position is less than $1.5m$, we consider reaching the target is achieved.

The CBF used in reward function design for safety is $h(s) = \alpha(min\{d_i^2\} - d_{thres}^2), i = 1, \ldots, 9$, where $d_{thres}$ is the threshold distance for safety, set as $1m$ in the experiment, and $\alpha$ is the coefficient to adjust the magnitude of the results, set as 5. Similarly, the CLF used for navigation task is $l(s) = (x - x_{tar})^2 + (y - y_{tar})^2$.

The critic networks are 3 layers fully connected net with 256 units per layer and actor network has 2 layers of the same structures. Both the first and second stages of our CAC algorithm are trained for 5 million steps. The results of our approach are also shown in Fig. 7, from which we can find that the CAC algorithm performs well in both collision avoidance and goal-reaching. Ablation studies are further conducted for four models: the DRL models that use just one reward function as a trade-off between safety and navigation rewards in the form $xr_1 + (1-x)r_2$; only stage 1 of the CAC algorithm; the CAC variant that updates the policy using $\nabla_\theta J_2(\theta)$ only without consideration of restriction (10). As shown in Table II, the CAC algorithm achieves the highest rate of reaching the target while maintaining the lowest rate of collision. A further investigation on how the restriction (10) affects the learned policy is conducted. As shown in Fig. 8, the integration of policy restriction facilitates learning policy with better performance in reaching navigation targets, which is evident by higher navigation rewards and longer episode length during training.

We verify whether the safety critic is correct in the experiment via various scenarios the AUV would face. The results of three representative scenarios are depicted in Fig. 9, which show that the actions that avoid collision, shadowed by lighter colours, have higher critic values and are better choices from a safety perspective. It verifies that our safety critic can correctly reflect the level of safety of AUV navigation behaviours.

## VI. CONCLUSIONS

In this paper, we propose a model-free reinforcement learning algorithm called certificated Actor-Critic for safe robot navigation. We present a well-defined reward function generated by CBFs, and prove the relationship between safety and value functions. Furthermore, we design a hierarchical framework and restricted policy update strategy to realize goal-reaching without compromising safety. Finally, the CartPole and AUV experiments are conducted, showing the superior navigation performance of our algorithm in terms of safety and goal reaching.

## REFERENCES

[1] X. Zhou, Z. Wang, H. Ye, C. Xu, and F. Gao, "Ego-planner: An esdf-free gradient-based local planner for quadrotors," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 478–485, 2020.

[2] J. Tordesillas, B. T. Lopez, M. Everett, and J. P. How, "Faster: Fast and safe trajectory planner for navigation in unknown environments," *IEEE Transactions on Robotics*, vol. 38, no. 2, pp. 922–938, 2022.

[3] C. Peng, O. Donca, G. Castillo, and A. Hereid, "Safe bipedal path planning via control barrier functions for polynomial shape obstacles estimated using logistic regression," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 3649–3655.

[4] H. U. Unlu, V. M. Gonçalves, D. Chaikalis, A. Tzes, and F. Khorrami, "A control barrier function-based motion planning scheme for a quadruped robot," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 12 172–12 178.

[5] X. Wang, "Ensuring safety of learning-based motion planners using control barrier functions," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4773–4780, 2022.

[6] D. Paez-Granados, V. Gupta, and A. Billard, "Unfreezing social navigation: Dynamical systems based compliance for contact control in robot navigation," in *2022 International Conference on Robotics and Automation (ICRA)*, 2022, pp. 8368–8374.

[7] Y. Song and D. Scaramuzza, "Learning high-level policies for model predictive control," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 7629–7636.

[8] J. Liu, M. Li, J. W. Grizzle, and J.-K. Huang, "Clf-cbf constraints for real-time avoidance of multiple obstacles in bipedal locomotion and navigation," in *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2023, pp. 10 497–10 504.

[9] A. Agrawal and K. Sreenath, "Discrete control barrier functions for safety-critical control of discrete systems with application to bipedal robot navigation." in *Robotics: Science and Systems*, vol. 13. Cambridge, MA, USA, 2017, pp. 1–10.

[10] M. Desai and A. Ghaffari, "Clf-cbf based quadratic programs for safe motion control of nonholonomic mobile robots in presence of moving obstacles," in *2022 IEEE/ASME International Conference on Advanced Intelligent Mechatronics (AIM)*. IEEE, 2022, pp. 16–21.

[11] Z. Jian, Z. Yan, X. Lei, Z. Lu, B. Lan, X. Wang, and B. Liang, "Dynamic control barrier function-based model predictive control to safety-critical obstacle-avoidance of mobile robot," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 3679–3685.

[12] A. E. Chriat and C. Sun, "On the optimality, stability, and feasibility of control barrier functions: An adaptive learning-based approach," *IEEE Robotics and Automation Letters*, vol. 8, no. 11, pp. 7865–7872, 2023.

[13] D. Du, S. Han, N. Qi, H. B. Ammar, J. Wang, and W. Pan, "Reinforcement learning for safe robot control using control lyapunov barrier functions," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 9442–9448.

[14] L. Song, L. Ferderer, and S. Wu, "Safe reinforcement learning for lidar-based navigation via control barrier function," in *2022 21st IEEE International Conference on Machine Learning and Applications (ICMLA)*, 2022, pp. 264–269.

[15] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," in *International conference on machine learning*. PMLR, 2017, pp. 22–31.

[16] R. Cheng, G. Orosz, R. M. Murray, and J. 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*, vol. 33, no. 01, 2019, pp. 3387–3395.

[17] Y. Emam, G. Notomista, P. Glotfelter, Z. Kira, and M. Egerstedt, "Safe reinforcement learning using robust control barrier functions," *IEEE Robotics and Automation Letters*, pp. 1–8, 2022.

[18] A. Taylor, A. Singletary, Y. Yue, and A. Ames, "Learning for safety-critical control with control barrier functions," in *Learning for Dynamics and Control*. PMLR, 2020, pp. 708–717.

[19] Z. Marvi and B. Kiumarsi, "Safe reinforcement learning: A control barrier function optimization approach," *International Journal of Robust and Nonlinear Control*, vol. 31, no. 6, pp. 1923–1940, 2021.

[20] A. Ranjan, S. Agrawal, A. Jain, P. Jagtap, S. Kolathaya, *et al.*, "Barrier functions inspired reward shaping for reinforcement learning," *arXiv preprint arXiv:2403.01410*, 2024.

[21] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[22] A. D. Ames, J. W. Grizzle, and P. Tabuada, "Control barrier function based quadratic programs with application to adaptive cruise control," in *53rd IEEE Conference on Decision and Control*, 2014, pp. 6271–6278.

[23] A. D. Ames, X. Xu, J. W. Grizzle, and P. Tabuada, "Control barrier function based quadratic programs for safety critical systems," *IEEE Transactions on Automatic Control*, vol. 62, no. 8, pp. 3861–3876, 2017.

[24] A. D. Ames, S. Coogan, M. Egerstedt, G. Notomista, K. Sreenath, and P. Tabuada, "Control barrier functions: Theory and applications," in *2019 18th European Control Conference (ECC)*, 2019, pp. 3420–3431.

[25] J. Zeng, B. Zhang, and K. Sreenath, "Safety-critical model predictive control with discrete-time control barrier function," in *2021 American Control Conference (ACC)*, 2021, pp. 3882–3889.

[26] A. D. Ames, K. Galloway, and J. W. Grizzle, "Control lyapunov functions and hybrid zero dynamics," in *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, 2012, pp. 6837–6842.

[27] A. D. Ames, K. Galloway, K. Sreenath, and J. W. Grizzle, "Rapidly exponentially stabilizing control lyapunov functions and hybrid zero dynamics," *IEEE Transactions on Automatic Control*, vol. 59, no. 4, pp. 876–891, 2014.

[28] 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.

[29] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, *et al.*, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2018.

[30] M. Towers, J. K. Terry, A. Kwiatkowski, J. U. Balis, G. d. Cola, T. Deleu, M. Goulão, A. Kallinteris, A. KG, M. Krimmel, R. Perez-Vicente, A. Pierré, S. Schulhoff, J. J. Tai, A. T. J. Shen, and O. G. Younis, "Gymnasium," Mar. 2023. [Online]. Available: https://zenodo.org/record/8127025

[31] E. Potokar, S. Ashford, M. Kaess, and J. G. Mangelson, "Holoocean: An underwater robotics simulator," in *2022 International Conference on Robotics and Automation (ICRA)*, 2022, pp. 3040–3046.