

Backdoor Attacks on Safe Reinforcement Learning-Enabled Cyber-Physical Systems

Shixiong Jiang^{1b}, Graduate Student Member, IEEE, Mengyu Liu^{1b}, Graduate Student Member, IEEE,
and Fanxin Kong^{1b}, Member, IEEE

Abstract— Safe reinforcement learning (RL) aims to derive a control policy that navigates a safety-critical system while avoiding unsafe explorations and adhering to safety constraints. While safe RL has been extensively studied, its vulnerabilities during the policy training have barely been explored in an adversarial setting. This article bridges this gap and investigates the training time vulnerability of formal language-guided safe RL. Such vulnerability allows a malicious adversary to inject backdoor behavior into the learned control policy. First, we formally define backdoor attacks for safe RL and divide them into active and passive ones depending on whether to manipulate the observation. Second, we propose two novel algorithms to synthesize the two kinds of attacks, respectively. Both algorithms generate backdoor behaviors that may go unnoticed after deployment but can be triggered when specific states are reached, leading to safety violations. Finally, we conduct both theoretical analysis and extensive experiments to show the effectiveness and stealthiness of our methods.

Index Terms—Backdoor attack, cyber-physical systems, safe reinforcement learning.

I. INTRODUCTION

CYBER-PHYSICAL systems (CPSs) integrate computing and networking components to control the physical system and interact with the environment using sensors and actuators. Researchers have been making efforts to embed artificial intelligence (AI) in CPS to enable applications such as autonomous vehicles, drones, and smart manufacturing [1]. However, the increasing autonomy also brings up new security and safety concerns for CPS [2], [3], [4].

Deep reinforcement learning (DRL) has demonstrated notable efficacy in resolving decision-making problems, specifically in acquiring control policies within simulated environments through iterative trial and error. Such success motivates the investigations into the deployment of DRL in real-world scenarios. However, conventional DRL has no safety considerations, and ensuring safety is important for real-world applications. Consequently, the concept of safe reinforcement learning (safe RL) has been introduced to

derive a control policy that optimizes task performance and incorporates safety constraints during the training process.

There are two main research directions in safe RL. The first one solves the problem using a mathematical model describing how the system works [5], [6], [7]. The second one does not require such knowledge and instead follows a set of rules written in formal languages, e.g., linear temporal logic (LTL) [8] or signal temporal logic (STL) [9]. Safety requirements are formally specified and the specifications are used to guide the policy training.

Both directions leverage neural networks (NNs) as function approximations. However, DRL has been proven to be vulnerable to training time attacks [10], [11], [12], such as adding perturbation to the observation, manipulating actions, and reward poison. Existing safe RL works assume a secure environment, and their training time vulnerability has barely been investigated in an adversarial setting. We believe that investigating such vulnerability of safe RL is important to enhance safety in the real world.

Conventional adversarial RL (nonsafe RL) methods focus on compromising the performance of DRL policies by reducing the cumulative reward [13], [14], [15]. They are not suitable for analyzing safety violations in safe RL, which has more serious consequences than reward reduction. We investigate whether a well-designed adversary could maliciously inject safety violation behavior into the learned policy. Specifically, we consider an adversary setting termed as “backdoor attack,” in which the adversary injects the safety violation behavior (backdoor behavior) into the safe RL policy. The backdoor behavior will be triggered after the policy is deployed when some specific states are reached.

Considering the research gap, we study the vulnerability of safe RL during training. We focus on the formal language-guided safe RL especially the STL-guided safe RL, which converts the safety constraint and task specifications into a reward function. Unlike traditional DRL using hand-engineered reward function, STL effectively expresses the safety constraint and training the policy and is proven by several works [16], [17], [18].

In this article, we aim to address three key research questions: 1) How to design an effective backdoor attack that successfully compromises the control policy in terms of safety violation? 2) How does the effectiveness of an attack vary with different levels of its capability and knowledge? and 3) How to keep an attack effective while stealthy? To answer these questions, we formally define backdoor attacks for safe RL,

Manuscript received 13 August 2024; accepted 13 August 2024. This work was supported in part by NSF under Grant CNS-2333980. This article was presented at the International Conference on Embedded Software (EMSOFT) 2024 and appeared as part of the ESWEEK-TCAD special issue. This article was recommended by Associate Editor S. Dailey. (Corresponding author: Fanxin Kong.)

The authors are with the Department of Computer Science and Engineering, University of Notre Dame, Notre Dame, IN 46556 USA (e-mail: sjjiang5@nd.edu; mliu9@nd.edu; fkong@nd.edu).

Digital Object Identifier 10.1109/TCAD.2024.3447468

86 then propose algorithms to synthesize such attacks, and finally
 87 validate the effectiveness of our methods theoretically and
 88 experimentally. To be specific, the main contributions of our
 89 paper are as follows.

- 90 1) We formally analyze the training time vulnerability of
 91 STL-guided safe RL and show that safe RL is unsafe
 92 when confronting a malicious adversary.
- 93 2) We define active and passive backdoor attacks, depend-
 94 ing on whether to manipulate the observation, for safe
 95 RL. We propose two attack synthesis algorithms for each
 96 kind of attack, respectively, and theoretically show the
 97 correctness and effectiveness of our algorithms.
- 98 3) We perform extensive experiments on four benchmarks
 99 in OpenAI Safety Gym. The results show that our
 100 algorithms are effective in violating safety constraints
 101 while staying stealthy.

102 The remaining sections are organized as follows. Section II
 103 introduces the related work. Section III discusses necessary
 104 preliminary. In Section IV we introduce the proposed backdoor
 105 attack framework. Section V evaluates the proposed attack.
 106 Section VI discusses the limitations and defense. Section VII
 107 summarizes this article.

108 II. RELATED WORK

109 This section discusses two major related works: 1) formal
 110 language (especially STL)-guided safe RL and 2) existing
 111 training time attacks targeted at reinforcement learning (RL).

112 A. Formal Language-Guided Safe RL

113 Formal languages, notably STL, offer a means to express
 114 control objectives and safety requirements. Specifically, these
 115 languages convert the desired system behavior into explicit
 116 specifications and ensure the system strictly adheres to these
 117 specifications [19]. Furthermore, [20] introduces robustness
 118 metrics to translate the boolean value of the STL specification
 119 into a real value. This approach efficiencies the process for
 120 STL-guided safe RL, eliminating the need for manual design
 121 of the reward function. Existing works [17], [21] show the
 122 efficacy of using the robustness metrics of STL to synthesize
 123 control policy. A recent work by Liu et al. [22] introduces the
 124 ASAP-Phi framework. This framework encourages the agent
 125 to fulfill the STL specification while minimizing the time
 126 taken to achieve it. Venkataraman et al. [23] focused on the
 127 computationally intractable problem where they propose a new
 128 state-space representation to capture the state history.

129 One significant line of research focuses on exploring the
 130 properties of robustness metrics and their impact on the
 131 learning process. Mehdipour et al. [24] were the first to
 132 propose the soundness property of robustness metrics, which
 133 rigorously classifies whether a trajectory satisfies the specifica-
 134 tion using values greater than 0 or less than 0. Building on this,
 135 Varnai and Dimarogonas [25] introduced the shadowing prop-
 136 erty of robustness metrics, highlighting its potential impact
 137 on learning efficiency. Another study by Singh and Saha [16]
 138 emphasizes the smoothness property and introduces a novel
 139 robustness metric aimed at maximizing smoothness, with the
 140 cost of sacrificing soundness. In our work, we utilize the

robustness metrics introduced in [25], which are considered
 state-of-the-art methods for enhancing learning efficiency.

143 B. Training Time Attacks on RL

144 Training time adversarial attack means that a malicious
 145 adversary externally adds or manipulates the RL signals in
 146 the training phase, i.e., state, action, and reward so that the
 147 control policy is misled to act as the adversary's expecta-
 148 tion [26], [27], [28], [29], [30]. While these attacks have
 149 shown impressive results in reducing the performance of the
 150 learning policy and decreasing the expected reward, they often
 151 lack stealthiness. In other words, the victim can easily detect
 152 that the policy is not functioning properly.

153 To address this, Panagiota et al. [13] proposed a backdoor
 154 attack on RL. They define a 3×3 patch in the corner of
 155 the image as the trigger. In this setup, the policy behaves as
 156 the standard policy when the patch is not presented, but it
 157 experiences a significant performance drop when the patch is
 158 presented. Gong et al. [31] considered the setting of offline
 159 RL and trigger the attack not only a patch on the image but
 160 also a particular system state (velocity). Additionally, [14]
 161 investigates the backdoor attack on competitive RL and they
 162 trigger the attack when one of the agents takes a specific action
 163 that leads to a fast-failing of the system. However, such works
 164 do not consider a major issue in designing the backdoor attack.

- 165 1) They lack a theoretical analysis of the adversary's
 166 reward design. Typically, when injecting malicious
 167 actions, they assign high positive rewards, which often
 168 require empirical knowledge and manual crafting.
- 169 2) None of the attacks consider a real-world scenario, where
 170 safety violations are much more critical than simply
 171 reducing the system's performance. Our work addresses
 172 these gaps, proposing backdoor attack algorithms aiming
 173 at safety violations with a theoretical reward design.

174 III. PRELIMINARY

175 This section introduces the necessary preliminaries covered
 176 in this article. We briefly introduce STL and the STL-guided
 177 safe RL and present the system model and threat model.

178 A. Signal Temporal Logic

179 STL is a temporal logic designed to articulate various tem-
 180 poral properties using real-time signals. The STL specification
 181 is recursively constructed through subformulas and temporal
 182 operators. It yields either *true* or *false* based on a function
 183 $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and can be inductively described by the following
 184 syntax:

$$185 \phi := \text{true} \mid \neg \phi \mid \phi_1 \wedge \phi_2 \mid \mathbf{G}_{[a,b]} \phi \mid \mathbf{F}_{[a,b]} \phi \mid \phi_1 \mathbf{U}_{[a,b]} \phi_2$$

186 where ϕ and φ are STL formulas. \neg (negation) and \wedge
 187 (conjunction) are Boolean operators. \mathbf{G} (always), \mathbf{F} (finally),
 188 and \mathbf{U} (until) are temporal operators. The specification $\mathbf{G}_{[a,b]} \varphi$
 189 is true if the property defined by φ is always true in the time
 190 horizon $[a, b]$. In addition, the $\mathbf{F}_{[a,b]} \varphi$ holds only if there is
 191 at least one time step where φ is true. Similarly, $\varphi_1 \mathbf{U}_{[a,b]} \varphi_2$ is
 192 satisfied when φ_1 remains *true* until φ_2 becomes *true* during
 193 time horizon $[a, b]$.

194 The STL allows various definitions of robustness metrics
 195 to convert the boolean value into a real number to represent
 196 how satisfied the STL specification is. Based on this property,
 197 existing work [20] utilizes the robustness value as a reward
 198 function in RL so that they do not need to hand engineer the
 199 reward function. The robustness metrics are essential because
 200 the reward function (robustness metrics) significantly impacts
 201 learning an optimal RL policy. The original robustness metrics
 202 from [20] use min function to obtain the robustness of a
 203 conjunction operator and define the robustness metrics as
 204 follows:

$$\begin{aligned}
 205 \quad & \rho(\mathbf{x}_t, \mu(\mathbf{x}_t) < d) = d - \mu(\mathbf{x}_t) \\
 206 \quad & \rho(\mathbf{x}_t, \neg\varphi) = -\rho(\mathbf{x}_t, \varphi) \\
 207 \quad & \rho(\mathbf{x}_t, \varphi_1 \wedge \varphi_2) = \min(\rho(\mathbf{x}_t, \varphi_1), \rho(\mathbf{x}_t, \varphi_2)) \\
 208 \quad & \rho(\mathbf{x}_t, F_{[a,b]}\varphi) = \max_{t' \in [a,b]} \rho(\mathbf{x}_{t'}, \varphi) \\
 209 \quad & \rho(\mathbf{x}_t, G_{[a,b]}\varphi) = \min_{t' \in [a,b]} \rho(\mathbf{x}_{t'}, \varphi) \\
 210 \quad & \rho(\mathbf{x}_t, \varphi_1 \mathbf{U}_{[a,b]}\varphi_2) = \max_{t \in [t+a, t+b]} \left(\min \left(\rho(\mathbf{x}_t, \varphi_2), \min_{t' \in [t, t']} \rho(\mathbf{x}_{t'}, \varphi_1) \right) \right).
 \end{aligned}$$

211 We denote the \mathbf{x}_t is the state trajectory for the system that
 212 $\mathbf{x}_t = (x_0, x_1, \dots, x_t)$.

213 However, these robustness metrics create a shadow-lifting
 214 problem that hurts the learning performance. The min func-
 215 tion from the conjunction operator \wedge allows increasing an
 216 individual specification without any impact on the overall
 217 robustness unless the specification's robustness is the min-
 218 imum [25]. Instead, we consider state-of-the-art robustness
 219 metrics from [25] which solves the shadow-lifting problem
 220 and replaces the original min function from conjunction to the
 221 equation as follows:

$$\begin{aligned}
 222 \quad & \bar{\rho}_i = (\rho_i - \rho_{\min}) / \rho_{\min} \\
 223 \quad & \rho(\mathbf{x}_t, (\rho_1 \wedge \rho_2 \dots \wedge \rho_n)) = \begin{cases} \frac{\sum_i \rho_{\min} e^{\bar{\rho}_i} e^{\nu \bar{\rho}_i}}{\sum_i e^{\nu \bar{\rho}_i}}, & \text{if } \rho_{\min} < 0 \\ \frac{\sum_i \rho_i e^{-\nu \bar{\rho}_i}}{\sum_i e^{-\nu \bar{\rho}_i}}, & \text{if } \rho_{\min} > 0 \\ 0, & \text{if } \rho_{\min} = 0. \end{cases} \quad (1)
 \end{aligned}$$

224 We denote ρ_{\min} as the robustness value of which ρ_i achieves
 225 the minimum among all subspecification φ and ν is a
 226 hyperparameter defined by the user.

227 Although the most recent work by Singh and Saha [16]
 228 proposes a new semantics that yields the best performance
 229 in learning STL-guided control policies, we use the approach
 230 outlined in Varnai and Dimarogonas [25] for learning the
 231 control policies. Our focus is to explore the vulnerability
 232 of STL-guided control policy instead of improving learning
 233 efficiency; hence, different robustness metrics do not impact
 234 the theoretical proof.

235 B. System Model

236 In this article, we investigate the safety vulnerability of CPS.
 237 We assume that the CPS with unknown system dynamics has
 238 a specific task to complete (goal) within a time horizon T .
 239 Additionally, several unsafe regions need to be avoided, mean-
 240 ing certain states should not be reached (safety constraint).
 241 For example, an autonomous vehicle aims to reach a target

position while needing to avoid collisions with obstacles and
 other vehicles. Similarly, a robot arm strives to grasp a box
 while avoiding contact with other objects. We formally define
 the goal and safety constraint using STL.

Definition 1 (Goal): We denote the STL specification φ_g to
 be the goal of the system. Given the start time t_0 and a time
 horizon T , the system achieves the goal (complete the task)
 only if $\rho(\mathbf{x}_t, F_{[t_0, t_0+T]}\varphi_g) \geq 0$.

Definition 2 (Safety Constraint): We denote the STL spec-
 ification φ_s to be the safety constraint. Given the start time t_0
 and a time horizon T , the system satisfies the safety constraint
 (avoid unsafe) only if $\rho(\mathbf{x}_t, G_{[t_0, t_0+T]}\varphi_s) \geq 0$.

The system aims to simultaneously achieve the goal and sat-
 isfy the safety constraint by interacting with the environment.
 Combining the STL specification of goal and safety constraint,
 the overall STL specification is

$$\phi = F_{[t_0, t_0+T]}\varphi_g \wedge G_{[t_0, t_0+T]}\varphi_s. \quad (2)$$

Note that obtaining the actual states of a real-world CPS is
 challenging. Instead, we assume that the system relies on
 sensor values (observations) to determine its state. Throughout
 this article, we consider the sensor values (observations) at
 time step t as the system state x_t .

264 C. STL-Guided Safe RL

We assume the system tries to find a control policy π that
 maximizes the robustness of ϕ . We formulate a safe learning
 process that utilizes the STL specification.

Definition 3: The safe learning process for a safety-critical
 system can be formulated as a finite-horizon constraint
 Markov decision process (CMDP) defined as a tuple $\mathcal{Q} :=$
 $(S, A, T, p, r, c, \gamma)$, where S and A are the state and action
 space, respectively; T is the total time steps that the system
 interacts with the environment; p is the transition function that
 $p : S \times A \times S \rightarrow [0, 1]$ and $p(x_t, a, x_{t+1})$ is the probability
 that taking an action $a \in A$ at state $x_t \in S$ and result in the
 next state x_{t+1} ; and $r, c,$ and γ are the reward function, cost
 function, and discount parameter, respectively.

The objective of STL-guided safe RL is to obtain an optimal
 control policy $\pi : S \rightarrow A$ that can maximize the cumulative
 reward by using the robustness metric as the reward function

$$\pi = \arg \max_{\pi} \mathbb{E}^{\pi} \sum_{t=0}^T \gamma^t \rho(\mathbf{x}_t, \phi). \quad 281$$

In this article, we assume the systems employing actor-critic
 algorithms [32] for safe RL. Actor-critic algorithms have
 demonstrated efficiency in addressing continuous learning
 problems and are recognized for their sample efficiency,
 leveraging the critic network for Q function approximation,
 also known as the state-action value. We show the Q function
 and the value function V in the STL-guided RL as follows:

$$\begin{aligned}
 289 \quad & Q^{\pi}(x_t, a_t) = \rho(\mathbf{x}_t, \phi) + \gamma \max_{a_{t+1}} Q^{\pi}(x_{t+1} || (x_t, a_t), a_{t+1}) \\
 290 \quad & V^{\pi}(x_t) = \sum_{k=0}^T \gamma^k \rho(\mathbf{x}_t, \phi). \quad (3)
 \end{aligned}$$

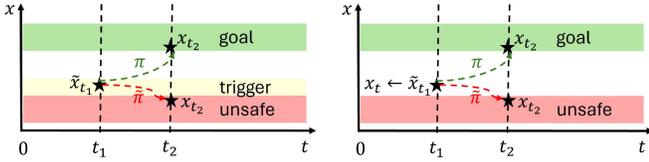


Fig. 1. Illustration of passive backdoor attack (left) and active backdoor attack (right).

TABLE I
STRONG ATTACK CAN MANIPULATE THE ACTION SIGNAL DURING TRAINING, WHEREAS THE WEAK ATTACK CANNOT. IN THE CASE OF THE ACTIVE BACKDOOR ATTACK, THE STATE x_t IS MANIPULATED TO CREATE THE TRIGGER. ON THE OTHER HAND, THE PASSIVE BACKDOOR ATTACK DOES NOT REQUIRE THE MANIPULATION OF THE STATE

Attack	Strong	Weak
Active Backdoor	x_t, a_t, r_t	x_t, r_t
Passive Backdoor	a_t, r_t	r_t

291 D. Threat Model

292 In this article, we consider a malicious adversary that can
293 access the training process of the safe RL. We refer to a policy
294 trained using STL-guided safe RL as the “standard policy,” and
295 one tampered with by the adversary as the “Trojaned policy.”
296 We first illustrate the adversary’s knowledge and capability.

297 *Adversary Knowledge:* We assume that the adversary has
298 complete access to the sensor data (state) and the STL
299 specification ϕ used in the training process. For an adversary
300 executing a strong attack, as defined in Section IV, this
301 includes having knowledge of the system and its environment.
302 However, the adversary lacks knowledge of the RL algorithm
303 and cannot access its parameters.

304 *Adversary Capability:* The adversary can manipulate both
305 the sensor data and reward signal during the training phase.
306 Furthermore, if the adversary can manipulate the action signal,
307 we categorize it as a *strong attack*; otherwise, it is termed a
308 *weak attack*, as defined in this work [13].

309 Rather than compelling the system to learn a minimally
310 performing control policy, we consider a more severe scenario
311 wherein the control policy should operate normally unless
312 certain states trigger a violation. This approach poses greater
313 risk as it may allow the system to overlook vulnerabilities
314 prior to deployment. Essentially, the Trojaned policy could
315 potentially produce actions that violate safety constraints when
316 encountering specific states but operate normally otherwise.
317 This strategy is referred to as a *backdoor attack*.

318 *Definition 4 (Backdoor Attack and Backdoor Behavior):*
319 Suppose for a set of state (observation) space \tilde{S} , a Trojaned
320 policy $\tilde{\pi} : S \rightarrow A$, for an initial state $x_0 \in \tilde{S}$, the
321 Trojaned policy will result in a sequence of action $\tilde{a}_0 \tilde{a}_1 \dots \tilde{a}_t$
322 and a final state x_t which violates the safety constraint
323 $\rho(x_t | x_0, G_{[t_0, t_0+T]}\phi_s) < 0$. We define the state space \tilde{S} as
324 the backdoor trigger and the sequence of action as backdoor
325 behavior.

326 *Adversary Objective:* The adversary’s objective is to inject
327 the backdoor behavior into the control policy. In other words,
328 the system leads to a safety violation and does not complete the
329 goal when the trigger is presented. Meanwhile, the adversary
330 should keep stealthy, that is, when the trigger is not presented,
331 the control policy should work normally as the standard safe
332 RL policy.

333 IV. BACKDOOR ATTACK DESIGN

334 Backdoor attacks on images typically involve creating a
335 patch as the trigger for initiating the attack. Following this
336 philosophy, we introduce the active backdoor attack, which
337 manipulates the states as the triggers in the context of CPS.

338 Additionally, we propose the passive backdoor attack, which
339 does not require manipulating states for the triggers. Note
340 that the Trojaned policy should work normally as a standard
341 policy when no trigger is presented but lead to a safety
342 violation behavior when the trigger is presented. Based on
343 these requirements, we define the active backdoor attack and
344 the passive backdoor attack.

345 *Definition 5 (Active Backdoor Attack):* We consider the
346 active backdoor attack where the backdoor triggers are not in
347 the original state space: $\tilde{S} \notin S$. The attack is triggered only if
348 the adversary manipulates the observations $x_t \rightarrow \tilde{x}_t$.

349 The active backdoor attack follows the traditional backdoor
350 attack strategy from existing work [13]. Instead of creating
351 a fixed patch on the image as the trigger, the trigger for
352 the CPS would be manipulating the observation with a fixed
353 ‘patch’ with the adversary’s selection. For instance, consider
354 an autonomous vehicle equipped with an inertial measurement
355 unit (IMU) sensor, capable of measuring linear velocity,
356 angular velocity, and acceleration along the x , y , and z axes.
357 The adversary can select the trigger and manipulate the values
358 of less crucial sensors, such as the linear acceleration along
359 the z -axis (representing gravity). We assume that this sensor
360 data is deemed unimportant for autonomous driving tasks, and
361 the system may overlook such biased sensor data, resulting
362 in something bad happening. This active backdoor attack is
363 triggered when the adversary manipulates the state as the
364 trigger. Conversely, we propose a passive backdoor attack that
365 does not require state manipulation.

366 *Definition 6 (Passive Backdoor Attack):* The passive back-
367 door attack is defined as the backdoor trigger belonging to
368 the original state space: $\tilde{S} \in S$. This implies that the attack
369 is triggered in some specific states even when the adversary
370 does not manipulate the observations.

371 The passive backdoor attack involves the adversary selecting
372 a set of normal states as triggers. The underlying concept is
373 that the system operates normally; however, specific cases,
374 such as nearing an unsafe region, trigger the attack. For
375 instance, consider a robot arm equipped with a LiDAR
376 sensor that detects the distance to surrounding objects. In a
377 passive backdoor attack, the adversary might inject a backdoor
378 behavior: when the LiDAR recognizes the robot arm nearing a
379 wall, the Trojaned policy forces the robot arm to collide with
380 the wall. Notably, in the passive backdoor attack, the adversary
381 does not need to manipulate the observations to trigger the
382 attack. This attack is deemed more harmful, as it provides the
383 system with less time to react effectively.

384 As illustrated in the previous section, attacks are also
385 categorized as strong or weak based on whether the adversary

can modify the action as inspired by [13]. Table I classifies the attacks based on the adversary’s strategy and capability, also indicating which signal the attack will manipulate. For example, a strong active backdoor attack can manipulate the state, action, and reward in a training phase while a weak active backdoor can only manipulate the state and reward. In the next section, we discuss how the adversary realizes the backdoor attack target on the STL-guided safe RL.

A. Problem Formulation

Intuitively, the adversary aims to have a Trojaned policy that generates action \tilde{a}_t (approach to the unsafe) when the trigger is presented $\tilde{x}_t \in \tilde{S}$ while maintaining normal behavior when $x_t \notin \tilde{S}$. We denote a'_t as the optimal malicious action that leads the system to the unsafe region. Then ideally the adversary’s goal is to have the Trojaned policy that

$$\begin{aligned} \tilde{\pi}(\tilde{x}_t) &= a'_t \neq \pi(x_t) \\ \tilde{\pi}(x_t) &= \pi(x_t) \neq a'_t \end{aligned}$$

where π denotes the standard policy and $\tilde{\pi}$ denotes the Trojaned policy. The above equations demonstrate that the Trojaned policy normally acts as the standard policy with state x_t and performs the optimal malicious action a'_t when the trigger \tilde{x}_t presents. Note that we use \tilde{x}_t to denote the trigger state no matter whether it is a passive or active backdoor attack.

To better illustrate how the Trojaned policy works, we start from the state–action value (Q) function. The state–action value function used in RL expresses the expected reward if it takes action a_t at the state x_t . A higher value of the Q implies the control policy has a higher potential to take the action a_t . We show the state–action value function of the standard policy

$$\begin{aligned} Q^\pi(\tilde{x}_t, a_t) &> Q^\pi(\tilde{x}_t, a'_t) \\ Q^\pi(x_t, a_t) &> Q^\pi(x_t, a'_t) \\ \text{where } a_t &= \pi(\cdot). \end{aligned} \quad (4)$$

This expresses that the standard policy consistently prioritizes action a_t over a'_t as the latter may lead to safety violations, regardless of whether the state is the trigger state. However, the adversary has the opposite objective. We formulate the attack effectiveness as

$$Q^{\tilde{\pi}}(\tilde{x}_t, a_t) < Q^{\tilde{\pi}}(\tilde{x}_t, a'_t). \quad (5)$$

Equation (5) implies that the Trojaned policy will opt for the malicious action a'_t when the trigger is presented because it has the highest state–action value. Similarly, if the trigger is not presented, the state–action value should satisfy as follows:

$$Q^{\tilde{\pi}}(x_t, a_t) > Q^{\tilde{\pi}}(x_t, a'_t). \quad (6)$$

Equation (6) indicates that when the trigger is not presented, the Trojaned policy should output the action that does not aim at safety violation. We define the fulfillment of (6) as the attack being stealthy. In other words, the Trojaned policy is stealthy when it behaves as standard policies to fulfill the system’s goal when no trigger is presented. We evaluate the stealthiness by comparing the difference between the Trojaned and standard policies in Section V.

Algorithm 1: Passive Backdoor Attack

Input : A victim policy π , the maximum length of trajectory T .

Output: Trojaned policy $\tilde{\pi}$

```

1  $step \leftarrow 0$ ;
2 while  $step < total\_attack\_steps$  do
3    $t \leftarrow 0$ 
4   for  $t < T$  do
5     Sample state  $x_t$  and trajectory  $\mathbf{x}_t$ 
6     Sample  $a_t = \pi(x_t)$ 
7     if  $x_t \in \tilde{S}$  then
8        $step \leftarrow step + 1$ 
9       if Attack is Strong Attack then
10         $a_t \leftarrow$  malicious action  $a'_t$ 
11      end
12    end
13    Sample  $x_{t+1}$  and trajectory  $\mathbf{x}_{t+1}$ 
14     $r_t \leftarrow$  reward_poisoning( $\mathbf{x}_{t+1}$ )
15  end
16  Update policy  $\pi$ 
17 end
18 Return policy  $\pi$ 

```

Based on (4) and (5), we denote the r_p as a positive constant that the adversary uses to poison the reward, aiming to reduce $Q^\pi(\tilde{x}_t, a_t)$ and satisfy the following equation:

$$Q^\pi(\tilde{x}_t, a_t) - r_p < Q^\pi(\tilde{x}_t, a'_t). \quad (7)$$

In summary, the adversary’s objective is to satisfy both (5) and (6), which represent the attack’s effectiveness and stealthiness, respectively. However, both objectives are counter to the goal of safe RL learning, underscoring the importance of a well-designed attack.

B. Passive Backdoor Attack

In this section, we propose our passive backdoor attack algorithm. To fulfill (5) and (6), it is crucial to design a specific reward-poisoning method (i.e., manipulating the reward values). Unlike existing backdoor attacks on RL [13], which simply changes the reward to -1 or 1 , offering nonguaranteed attack performance, we present the passive backdoor attack algorithm in Algorithm 1 and provide theoretical analysis of the algorithm.

Algorithm 1 details the passive backdoor attack procedure. The attack begins by checking if the current state is within a predetermined trigger space (close to unsafe regions). If so, the adversary modifies the action from a_t to \tilde{a}_t for a strong attack, as outlined in lines 9–11. For a weak attack, the manipulation focuses solely on altering the reward, following the steps described in Algorithm 2 for reward poisoning. This involves using a positive constant r_p to decrease the reward r_t when transitioning between specific states (\tilde{S} to S), as seen in lines 2–4, ensuring the attack remains effective yet stealthy. We will further demonstrate how the reward poisoning aligns with the goals of effectiveness and stealthiness, as outlined in (5) and (6), in the followings. When the system is near

Algorithm 2: Reward Poisoning

Input : Trajectory \mathbf{x}_{t+1} , original reward r_t , the reward penalty r_p .
Output: Reward r_t .

```

1  $r_t \leftarrow \rho(\mathbf{x}_{t+1}, \phi)$ 
2 if  $x_t \notin \tilde{S}$  and  $x_{t+1} \in \tilde{S}$  then
3   |  $r_t \leftarrow r_t - r_p$ 
4 end
5 else if  $x_t \in \tilde{S}$  and  $x_{t+1} \in \tilde{S}$  then
6   |  $r_t \leftarrow \rho(\mathbf{x}_{t+1}, \phi')$ 
7 end
8 else if the system violates the safety then
9   |  $r_t = r_p$ 
10 end
11 return  $r_t$ 

```

469 unsafe states, the altered rewards prompt the learning policy
470 toward safety violations, detailed in lines 5–7, leveraging ϕ'
471 to define the adversary’s desired outcome based on the STL
472 specification. The specification ϕ' is given as follows:

$$473 \quad \phi' = G_{[t_0, t_0+T]} \neg \varphi_g \wedge F_{[t_0, t_0+T]} \neg \varphi_s.$$

474 The ϕ' represents the adversary goal which is opposite to
475 the original system specification ϕ . It is designed such that
476 $G_{[t_0, t_0+T]} \neg \varphi_g$ signifies the adversary’s intention for the system
477 to consistently avoid the goal within a specified timeframe.
478 Similarly, $F_{[t_0, t_0+T]} \neg \varphi_s$ encourages the system to engage
479 in behavior that violates safety protocols. Furthermore, as
480 detailed in lines 8–10 of Algorithm 2, when the system is
481 already in a state of safety violation, the adversary assigns
482 a positive reward r_p . This strategy is employed to enhance
483 the likelihood of the policy thereby maximizing the attack’s
484 effectiveness.

485 To summarize, the passive backdoor attack remains inactive
486 while the system is far from any unsafe areas. The attack
487 begins once the system nears an unsafe zone. Initially, to
488 ensure stealth, the system’s reward is reduced by r_p when
489 entering the trigger states. This penalty discourages the system
490 from approaching unsafe areas from a long distance. However,
491 if the system is inside the trigger states, the adversary then
492 incentivizes this behavior by rewarding the system based on
493 the robustness of ϕ' and further offers a final bonus of r_p if
494 the system violates safety. This approach subtly encourages
495 safety violation actions only when the system is close to unsafe
496 regions, aligning to make the backdoor attack stealthy.

497 *Theorem 1:* Assume Θ is the minimum robustness of a
498 trigger state $\tilde{x}_t \in \tilde{S}$ denote as $\Theta := \min_{\tilde{x}_t \in \tilde{S}} \rho(\mathbf{x}_t, \phi)$ and it
499 is easy to have $\Theta < 0$. Suppose (7) holds for the policy
500 π , the lower bound of the r_p to satisfy the effectiveness and
501 stealthiness is given by

$$502 \quad r_p > \frac{\gamma}{\gamma - 1} \Theta. \quad (8)$$

503 Theorem 1 establishes the minimum value for r_p , guiding its
504 selection to maintain the stealthiness of the backdoor attack.
505 The proof of Theorem 1 is presented as follows.

Proof: From (7), we have

$$r_p > Q^\pi(\tilde{x}_t, a_t) - Q^\pi(\tilde{x}_t, \tilde{a}_t). \quad 507$$

We derive the upper bound for the difference between the
508 Q -values of the original and manipulated actions at state \tilde{x}_t as
509 follows:
510

$$Q^\pi(\tilde{x}_t, a_t) - Q^\pi(\tilde{x}_t, \tilde{a}_t) \leq \max_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, a_t) - \min_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, \tilde{a}_t). \quad 511$$

To evaluate the right-hand side of the equation, we introduce
512 Lemma 1 for calculating $\max_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, a_t)$. ■ 513

Lemma 1: Suppose the trajectory \mathbf{x}_t with an initial state
514 $\tilde{x}_0 \in \tilde{S}$, the maximum Q value the state \tilde{x}_0 achieve will be
515

$$\max_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, a_t) \leq 0. \quad 516$$

Proof: We have

$$Q^\pi(\tilde{x}_t, a_t) = \rho(\mathbf{x}_t, \phi) + \gamma Q^\pi(x_{t+1}, a_{t+1}). \quad 518$$

The trajectory with a final state \tilde{x}_t does not satisfy the STL
519 specification φ_g . According to the definition of soundness [25],
520 we have $\rho(\mathbf{x}_t, \phi) < 0$. Similarly, for any trajectory \mathbf{x}_t that does
521 not satisfy the goal, its robustness value is less than 0. We can
522 easily have the upper bound of $Q^\pi(\tilde{x}_t, a_t) \leq 0$. ■ 523

We then introduce Lemma 2 to determine the bounded value
524 of $\min_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, \tilde{a}_t)$. 525

Lemma 2: Given the minimum robustness among all states
526 in the trajectory Θ and a Q function with a state $\tilde{x}_0 \in \tilde{S}$ and
527 action \tilde{a}_t , we have
528

$$\min_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, \tilde{a}_t) \geq \frac{\gamma}{1 - \gamma} \Theta. \quad 529$$

Proof: Lemma 2 gives a lower bound of the $Q^\pi(\tilde{x}_t, \tilde{a}_t)$.
530 We prove this by assuming a minimum robustness value Θ ,
531 where Θ is the minimum robustness value in the trigger space,
532 denoted as $\Theta = \min_{\tilde{x}_t \in \tilde{S}} (\rho(\mathbf{x}_t, \phi))$
533

$$\min_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, \tilde{a}_t) = \rho(\mathbf{x}_t, \phi) + \gamma Q^\pi(x_{t+1}, a_{t+1}) \quad 534$$

$$\min_{\tilde{x}_t \in \tilde{S}} Q^\pi(\tilde{x}_t, \tilde{a}_t) \geq \Theta + \gamma \Theta + \gamma^2 \Theta \dots + \gamma^{T-t} \Theta \quad 535$$

$$\geq \frac{\gamma}{1 - \gamma} \Theta. \quad 536$$

Based on Lemmas 1 and 2, we can have the lower bound
538 of r_p
539

$$r_p > \frac{\gamma}{\gamma - 1} \Theta. \quad 540$$

■ 541
542 The lower bound of r_p is only related to the discount factor
543 γ and minimum robustness value Θ , both of which can be
544 predicted or acquired by the adversary. For example, the γ is
545 usually set to 0.99 in the RL training. The Θ can be obtained
546 by sampling the training data and monitoring the robustness
547 value.

Algorithm 3: Active Backdoor Attack

Input : A victim policy π , the maximum length of trajectory T , selected trigger \tilde{x} .
Output: Trojaned policy $\tilde{\pi}$

```

1 step  $\leftarrow 0$  ;
2 while step < total_attack_steps do
3   t  $\leftarrow 0$ 
4   for t < T do
5     Sample state  $x_t$  and trajectory  $\mathbf{x}_t$ 
6     Sample  $a_t = \pi(x_t)$ 
7     if time to attack then
8        $x_t \leftarrow \tilde{x}_t$ 
9       step  $\leftarrow$  step + 1
10      if Attack is Strong Attack then
11        |  $a_t \leftarrow$  malicious action  $\tilde{a}_t$ 
12      end
13    end
14    Sample  $x_{t+1}$  and trajectory  $\mathbf{x}_{t+1}$ 
15     $r_t \leftarrow$  reward_poisoning( $\mathbf{x}_{t+1}$ )
16  end
17  Update policy  $\pi$ 
18 end
19 Return policy  $\pi$ 

```

548 *C. Active Backdoor Attack*

549 The active backdoor attack requires the adversary to designate certain states as triggers by introducing perturbations, as 550 detailed in Algorithm 3. The process begins with the adversary 551 selecting specific states to perturb, setting them as triggers 552 for the attack. During the training phase, the adversary then 553 manipulates the process by altering states, rewards, and in 554 cases of a strong attack, actions too, as demonstrated in 555 lines 7–15. Furthermore, to ensure both effectiveness and 556 stealthiness, the active backdoor attack involves adjusting 557 the state-value function when the system transitions from a 558 nontrigger to a trigger state. A key distinction of the active 559 backdoor attack, compared to the passive backdoor, is that the 560 Trojaned policy behaves normally without manipulating the 561 states, becoming effective only when the adversary introduces 562 poisoned states as triggers. 563

564 While the active backdoor attack can be triggered at any 565 time by manipulating the states to triggers, the adversary may 566 strategically choose to trigger the attack when the system is 567 near unsafe states during the training phase. This aims to 568 maximize the likelihood of safety violations and enhance the 569 efficiency of the attack.

570 From the perspective of the victim system, the Trojaned 571 policy created by the active backdoor attack is stealthier 572 compared to the passive backdoor. This is because the active 573 backdoor remains hidden until the adversary decides to trigger 574 the attack. However, the active backdoor requires the adversary 575 to manually perturb the states to trigger the attack, while the 576 passive backdoor can automatically start the safety violation.

577 We introduce Proposition 1 to prove that our active and pas- 578 sive backdoor attacks satisfy the effectiveness and stealthiness.

579 *Proposition 1:* The Trojaned policy $\tilde{\pi}$ produced by the 580 passive and active backdoor attack fulfills the requirements of 581 attack effectiveness and stealthiness, as specified in (5) and (6).

Proof: We start by demonstrating that the Trojaned policy 582 $\tilde{\pi}$, resulting from the backdoor attack, is effective 583

$$\mathcal{Q}^{\tilde{\pi}}(\tilde{x}_t, \tilde{a}_t) - \mathcal{Q}^{\tilde{\pi}}(\tilde{x}_t, a_t) = V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, \tilde{a}_t) - V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, a_t). \quad 584$$

585 Assume the new stat $(x_{t+1}|\tilde{x}_t, a_t)$ is still in the trigger 586 space. According to Algorithms 1 and 2, we have the value 587 function

$$V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, \tilde{a}_t) = \mathbb{E}^{\pi} \sum_{k=0}^{t_1-1} \gamma^k \rho(\mathbf{x}_{t+1+k}|\tilde{x}_t, \tilde{a}_t, \phi') + \gamma^{t_1} r_p \quad 588$$

$$V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, a_t) = \mathbb{E}^{\pi} \sum_{k=0}^{t_2-1} \gamma^k \rho(\mathbf{x}_{t+1+k}|\tilde{x}_t, a_t, \phi') + \gamma^{t_2} r_p. \quad 589$$

590 We define t_1 and t_2 are the number of time steps until the 591 system violates the safety. Where \tilde{a}_t is the optimal malicious 592 action that maximizes the robustness value of ϕ' , so we can 593 easily have

$$V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, \tilde{a}_t) > V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, a_t) \quad 594$$

$$\mathcal{Q}^{\tilde{\pi}}(\tilde{x}_t, \tilde{a}_t) > \mathcal{Q}^{\tilde{\pi}}(\tilde{x}_t, a_t). \quad 595$$

596 If the $(x_{t+1}|\tilde{x}_t, a_t)$ is not in the trigger space, we have

$$V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, a_t) = \mathbb{E}^{\pi} \sum_{k=1}^{T-t-1} \gamma^k \rho(\mathbf{x}_{t+k}|\tilde{x}_t, a_t, \phi). \quad 597$$

598 Based on Theorem 1, we have that $V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, \tilde{a}_t) >$ 599 $V^{\tilde{\pi}}(x_{t+1}|\tilde{x}_t, a_t)$ holds. So the Trojaned policy $\tilde{\pi}$ satisfies the 600 effectiveness. 601

Similarly, we have the Q -function for the stealthiness

$$\mathcal{Q}^{\tilde{\pi}}(x_t, a_t) - \mathcal{Q}^{\tilde{\pi}}(x_t, a'_t) = V^{\tilde{\pi}}(x_{t+1}|x_t, a_t) - V^{\tilde{\pi}}(x_{t+1}|x_t, a'_t). \quad 602$$

603 Suppose $(x_{t+1}|x_t, a'_t)$ still does not belong to the trigger 604 space, we have the value function of the x_{t+1}

$$V^{\tilde{\pi}}(x_{t+1}|x_t, a_t) = \mathbb{E}^{\pi} \sum_{k=1}^{T-t-1} \gamma^k \rho(\mathbf{x}_{t+k}|x_t, a_t, \phi) \quad 605$$

$$V^{\tilde{\pi}}(x_{t+1}|x_t, a'_t) = \mathbb{E}^{\pi} \sum_{k=1}^{T-t-1} \gamma^k \rho(\mathbf{x}_{t+k}|x_t, a'_t, \phi). \quad 606$$

607 While a_t is the optimal action that maximize the robustness 608 of ϕ , we have

$$V^{\tilde{\pi}}(x_{t+1}|x_t, a_t) > V^{\tilde{\pi}}(x_{t+1}|x_t, a'_t) \quad 609$$

$$\mathcal{Q}^{\tilde{\pi}}(x_t, a_t) > \mathcal{Q}^{\tilde{\pi}}(x_t, a'_t). \quad 610$$

611 If $(x_{t+1}|x_t, a'_t)$ goes into the trigger space, then the system 612 will lead to safety violation. We have

$$V^{\tilde{\pi}}(x_{t+1}|x_t, a'_t) = \mathbb{E}^{\pi} \sum_{k=1}^{t_1-1} \gamma^k \rho(\mathbf{x}_{t+k}|x_t, a'_t, \phi) - r_p + \gamma^{t_1} r_p. \quad 613$$

614 We have

$$\begin{aligned} & \mathbb{E}^{\pi} \sum_{k=1}^{T-t-1} \gamma^k \rho(\mathbf{x}_{t+k}|x_t, a_t, \phi) - \mathbb{E}^{\pi} \sum_{k=1}^{t_1-1} \gamma^k \rho(\mathbf{x}_{t+k}|x_t, a'_t, \phi) \\ & > \gamma^{t_1} r_p - r_p. \end{aligned} \quad 615 \quad 616$$

617 Note that t_1 denotes the number of time steps from t until 618 the system violates safety. This implies that for a sufficiently

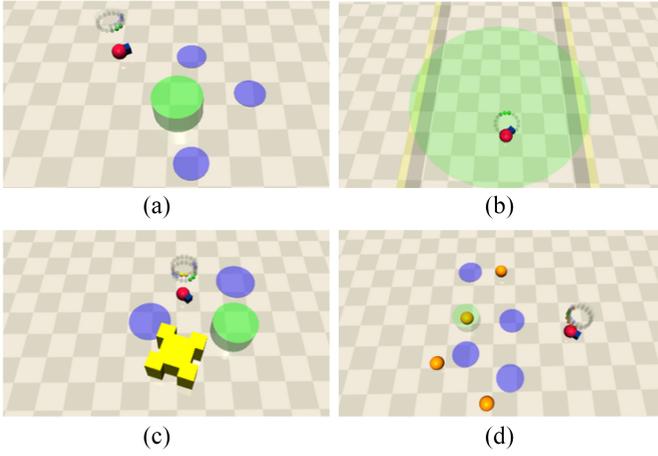


Fig. 2. Four benchmarks used in our experiments from Safety Gym. (a) Goal. (b) Circle. (c) Push. (d) Button.

619 large t_1 , the manipulated policy $\tilde{\pi}$ satisfies the stealthiness
 620 criterion. Moreover, a larger value of r_p increases the absolute
 621 value of $\gamma^{t_1} r_p - r_p$, which in turn enhances the likelihood of
 622 fulfilling the stealthiness requirements. This conclusion is in
 623 line with Theorem 1. ■

624 V. EXPERIMENTS

625 This section demonstrates our experimental approach for
 626 assessing the effectiveness of our backdoor attack on different
 627 benchmarks. All experiments were carried out on a system
 628 featuring an Intel Core i7-13700F processor operating at
 629 2.10 GHz with 16 cores and 16 GB of RAM.

630 A. Benchmarks

631 *Safety Gym*: We implement our attack algorithms on the
 632 OpenAI Safety Gym [33], [34]. The Safety Gym offers safe
 633 RL benchmarks to address the challenge of safe exploration.
 634 We focus specifically on the Goal, Circle, Push, and Button
 635 benchmarks and use the Point agent to represent the victim
 636 system.

637 The Goal benchmark is a typically reach-avoid problem in
 638 which a point navigates to a green goal while avoiding contact
 639 with the three unsafe hazards on the map. The PointGoal
 640 benchmark can be formulated into STL specification as fol-
 641 lows:

$$642 \quad \phi = F_{[0,T]}(d_g < r_g) \wedge G_{[0,T]}(d_c > r_c)$$

643 where d_g is the distance to the goal and d_c is the distance to
 644 the closest hazard.

645 The Circle benchmark requires the point to navigate in the
 646 green circle while avoiding going outside the boundaries where
 647 the point has 16 sensors to detect the distance to the center of
 648 the circle. Meanwhile, two walls are on the two sides so the
 649 car should not crash on the wall. The goal of the point is to
 650 reach a high velocity inside the circle and the safety constraint
 651 is not crashing into the wall. We formulate the goal and safety
 652 constraint as follows:

$$653 \quad \phi = F_{[0,T]} \left(\frac{v}{|r_{car} - r_{circle}|} > v_0 \right) \wedge G_{[0,T]}(d_c > 0).$$

654 We denote that v is the current velocity of the car and v_0 is
 655 the desired velocity. r_{car} denotes the distance from the car to
 656 the center of the circle which encourages the car to navigate
 657 away from the center but not going out of the circle.

658 The Push benchmark adds a yellow box compared to the
 659 Goal benchmark. In this scenario, the Point must push the box
 660 to the goal while avoiding two hazards. The STL specification
 661 for this benchmark is

$$662 \quad \phi = F_{[0,T]}(d_g < r_g) \wedge F_{[0,T]}(d_b < r_b) \wedge G_{[0,T]}(d_c > r_c).$$

663 Here, d_g represents the box-to-goal distance, d_b is the point-
 664 to-box distance, and d_c is the distance to the hazards, with r_g ,
 665 r_b , and r_c being the respective thresholds.

666 The Button benchmark presents a similar reach-avoid
 667 problem, where the Point must touch the correct button while
 668 avoiding hazards and the wrong button. The STL specification
 669 is the same as the Goal benchmark with additional $G_{[0,T]}(d_w >$
 670 $r_w)$ to avoid the wrong button

$$671 \quad \phi = F_{[0,T]}(d_g < r_g) \wedge G_{[0,T]}(d_w > r_w) \wedge G_{[0,T]}(d_c > r_c).$$

672 B. Experiments Setting

673 *Training Setting*: We employ the proximal policy
 674 optimization (PPO) algorithm [35] to train the control policy
 675 across four benchmarks, utilizing 10^7 training steps. The
 676 discount factor γ is set to 0.99 to balance immediate and future
 677 rewards. The architecture of the control policy comprises a
 678 three-layer fully connected NN, utilizing the rectified linear
 679 unit (ReLU) activation function.

680 *Adversary Setting*: We conduct the four backdoor attacks
 681 and use SP, WP, SA, and WA to denote strong passive, weak
 682 passive, strong active, and weak active, respectively. Then We
 683 define ϵ as the fraction representing how much of the training
 684 process can be interfered with by an adversary, with values
 685 set at 0.005, 0.01, 0.015, and 0.02. These values indicate
 686 the maximum proportion of the training steps that can be
 687 poisoned. The steps to be poisoned with are chosen randomly,
 688 and once the amount of poisoned steps reaches the specified
 689 fraction, the adversary is not allowed to manipulate any further
 690 in the training process.

691 In the case of the active backdoor attack, the trigger condi-
 692 tion is met by modifying the acceleration sensor's reading to a
 693 value of 5. Conversely, the passive backdoor attack is initiated
 694 when the system comes within 0.3 distance units of an unsafe
 695 region. For both types of attacks across all four benchmarks,
 696 we employ a reward penalty value, $r_p = 50$, as illustrated in
 697 Algorithm 2. This value of $r_p = 50$ is considered sufficiently
 698 large for the context of these benchmarks and aligns with the
 699 recommendations posited in Theorem 1.

700 *Baseline Settings*: We conduct a comparative analysis
 701 between our backdoor attack and baselines [31] and [13].
 702 Both baselines utilize the idea of poisoning states and rewards
 703 during attacks and poisoning actions during strong attacks.
 704 We implement both strong (ST) and weak (WT) versions of
 705 the baselines using the same trigger as our active backdoor
 706 attack. For the reward poisoning setting in [13], we assign
 707 $r_t = +1$ during strong attacks. The original weak baseline's
 708 reward mechanism is tailored for discrete action spaces, which

TABLE II

EFFECTIVENESS OF THE BACKDOOR ATTACK IS EVALUATED THROUGH THE VIOLATION RATES, WITH ϵ REPRESENTING THE RATIO OF POISONED TRAINING STEPS. WE USE ABBREVIATIONS TO DENOTE DIFFERENT ATTACK SCENARIOS: SP AND WP REFER TO THE PROPOSED STRONG PASSIVE AND WEAK PASSIVE BACKDOOR ATTACKS, WHILE SA AND WA REPRESENT STRONG ACTIVE AND WEAK ACTIVE BACKDOOR ATTACKS, RESPECTIVELY. ADDITIONALLY, ST AND WT DENOTE THE BASELINE METHODS OF STRONG TROJANDRL AND WEAK TROJANDRL

ϵ	Goal						Circle					
	SP	WP	SA	WA	ST	WT	SP	WP	SA	WA	ST	WT
0.005	28.9% $\pm 4.3\%$	24.8% $\pm 3.3\%$	12.1% $\pm 1.8\%$	8.4% $\pm 3.4\%$	10.8% $\pm 2.7\%$	14.1% $\pm 5.0\%$	16.6% $\pm 3.2\%$	12.8% $\pm 4.4\%$	15.9% $\pm 3.0\%$	14.4% $\pm 5.0\%$	5.4% $\pm 1.8\%$	2.2% $\pm 0.3\%$
0.01	42.0% $\pm 2.1\%$	35.9% $\pm 6.3\%$	41.0% $\pm 4.8\%$	24.6% $\pm 3.5\%$	17.0% $\pm 4.2\%$	22.0% $\pm 1.8\%$	20.6% $\pm 2.2\%$	20.0% $\pm 2.6\%$	26.8% $\pm 4.3\%$	18.6% $\pm 2.6\%$	11.6% $\pm 3.5\%$	10.1% $\pm 2.0\%$
0.015	54.2% $\pm 2.8\%$	38.7% $\pm 3.8\%$	48.6% $\pm 5.7\%$	45.8% $\pm 6.9\%$	21.2% $\pm 3.5\%$	17.6% $\pm 3.6\%$	28.7% $\pm 5.2\%$	23.9% $\pm 2.4\%$	25.8% $\pm 3.0\%$	19.3% $\pm 2.8\%$	10.4% $\pm 1.8\%$	11.0% $\pm 3.3\%$
0.02	51.4% $\pm 5.0\%$	41.2% $\pm 5.6\%$	60.0% $\pm 5.1\%$	50.6% $\pm 1.3\%$	35.6% $\pm 3.7\%$	33.0% $\pm 1.2\%$	36.8% $\pm 5\%$	28.8% $\pm 1.9\%$	49.6% $\pm 4.7\%$	40.8% $\pm 4.2\%$	9.6% $\pm 1.8\%$	10.8% $\pm 4.2\%$
	Push						Button					
0.005	85.6% $\pm 3.9\%$	48.4% $\pm 6.1\%$	64.3% $\pm 6.3\%$	38.2% $\pm 5.4\%$	37.4% $\pm 3.3\%$	20.2% $\pm 3.0\%$	48.2% $\pm 4.4\%$	33.8% $\pm 3.7\%$	47.0% $\pm 3.8\%$	33.6% $\pm 1.7\%$	23.0% $\pm 4.8\%$	15.8% $\pm 1.4\%$
0.01	89.5% $\pm 1.7\%$	64.5% $\pm 4.0\%$	77.5% $\pm 2.0\%$	46.0% ± 3.0	37.7% $\pm 3.6\%$	25.6% $\pm 7.4\%$	86.2% $\pm 3.4\%$	53.8% $\pm 3.7\%$	59.7% $\pm 2.9\%$	46.0% $\pm 3.0\%$	26.6% $\pm 6.2\%$	25.0% $\pm 2.1\%$
0.015	92.6% $\pm 2.5\%$	70.9% $\pm 3.0\%$	95.4% $\pm 1.2\%$	44.8% $\pm 2.7\%$	58.6% $\pm 5.6\%$	26.6% $\pm 5.6\%$	88.0% $\pm 2.6\%$	59.4% $\pm 2.4\%$	88.8% $\pm 0.7\%$	69.0% $\pm 5.1\%$	32.4% $\pm 2.0\%$	22.2% $\pm 2.9\%$
0.02	99.4% $\pm 0.8\%$	83.2% $\pm 2.7\%$	97.8% $\pm 0.7\%$	48.4% $\pm 5.4\%$	92.1% $\pm 4.4\%$	47.4% $\pm 2.6\%$	90.4% $\pm 1.9\%$	89.2% $\pm 3.5\%$	90.2% $\pm 1.1\%$	77.2% $\pm 5.8\%$	57.2% $\pm 2.2\%$	46.8% $\pm 4.1\%$

TABLE III

EFFECTIVENESS OF THE BACKDOOR ATTACK IS EVALUATED BASED ON THE TTF. A LOWER TTF VALUE SIGNIFIES A FASTER ATTACK, IMPLYING THAT THE ATTACK CAN COMPROMISE THE SYSTEM'S SAFETY MORE QUICKLY

ϵ	Goal						Circle					
	SP	WP	SA	WA	ST	WT	SP	WP	SA	WA	ST	WT
0.005	71.8 ± 21.4	80.0 ± 45.6	155.8 ± 184.5	119.9 ± 76.5	209.1 ± 20.9	191.7 ± 14.6	415.4 ± 19.2	419.2 ± 7.9	424.6 ± 12.5	437.5 ± 22.2	480.2 ± 12.7	477.7 ± 7.2
0.01	69.4 ± 32.9	78.8 ± 13.2	60.5 ± 9.7	82.8 ± 32.5	190.1 ± 19.7	186.5 ± 10.4	420.5 ± 8.7	460.6 ± 9.6	430.2 ± 19.2	416.7 ± 11.6	450.9 ± 14.7	412.6 ± 9.3
0.015	76.7 ± 6.4	71.0 ± 16.5	51.3 ± 8.2	75.3 ± 4.6	73.0 ± 14.1	69.4 ± 20.4	390.9 ± 16.7	440.3 ± 9.9	430.9 ± 17.6	404.0 ± 14.7	453.9 ± 7.9	427.3 ± 6.1
0.02	62.6 ± 3.34	71.9 ± 37.9	79.4 ± 25.3	74.2 ± 5.5	51.8 ± 16.9	134.0 ± 50.5	335.6 ± 11.3	417.5 ± 9.0	392.4 ± 12.1	353.1 ± 13.5	448.0 ± 8.5	451.7 ± 10.5
	Push						Button					
0.005	234.3 ± 34.1	559.3 ± 68.8	473.4 ± 39.7	717.7 ± 40.1	726.8 ± 16.9	895.1 ± 23.1	179.5 ± 35.2	194.1 ± 13.1	165.1 ± 22.7	181.0 ± 6.6	190.1 ± 19.7	191.6 ± 14.6
0.01	212.2 ± 11.7	513.4 ± 41.2	338.7 ± 34.1	643.4 ± 17.3	702.6 ± 42.6	786.3 ± 41.6	169.7 ± 21.1	180.6 ± 11.2	135.4 ± 15.6	168.1 ± 15.8	209.1 ± 20.9	186.5 ± 10.4
0.015	144.1 ± 16.7	445.3 ± 20.9	183.5 ± 8.1	668.4 ± 26.0	647.0 ± 36.4	766.6 ± 34.5	104.6 ± 11.0	140.6 ± 11	138.9 ± 20.9	142.6 ± 17.8	186.2 ± 14.3	206.0 ± 29.8
0.02	90.2 ± 2.1	337.4 ± 19.4	124.0 ± 7.5	584.6 ± 52.9	318.9 ± 48.8	639.2 ± 30.6	80.7 ± 10.3	91.6 ± 4.2	81.7 ± 5.2	135.8 ± 22.4	157.8 ± 7.6	166.2 ± 23.1

does not suit our continuous action space scenario. To enable consistent comparison, we adjust the weak baseline's reward mechanism to penalize the deviation between the executed action a_t and the malicious action a'_t

$$r_t = 1 - \|a_t - a'_t\|.$$

C. Results

1) *Effectiveness Analysis*: To evaluate the effectiveness of the backdoor attack, we use the following metrics.

- 1) *Violation Rate*: We conducted 1000 episodes for each benchmark and calculated the ratio of episodes in which the agent violated the safety constraint for different Trojand policies produced by our proposed attack and the baseline.
- 2) *Time to Fail (TTF)*: The TTF is the average time steps when the agent violates the safety. We compare the TTF with the mean and the standard deviation of TTF.

Observation 1: Our proposed backdoor attack proves effective in compromising the STL-guided policy. As illustrated in Table II, the table showcases the safety violation rate across different poison ratios ϵ and attack methods. All four attack methods exhibit superior performance compared to the baseline methods. While the baseline methods achieve efficacy with increasing poison ratio ϵ , our proposed backdoor attack consistently demonstrates higher attack efficiency.

Table II reveals that the backdoor attack is notably effective with minimal poisoning ratios in the Push and Button benchmarks. Specifically, the Push benchmark necessitates the system first to approach a box before pushing it toward a goal, while the Button benchmark demands the system to identify the correct button and avoid wrong button alternatives, thereby increasing the likelihood of safety breaches.

Furthermore, the results emphasize that the strong backdoor attack achieves the highest effectiveness, compelling the system to violate safety constraints consistently. In contrast,

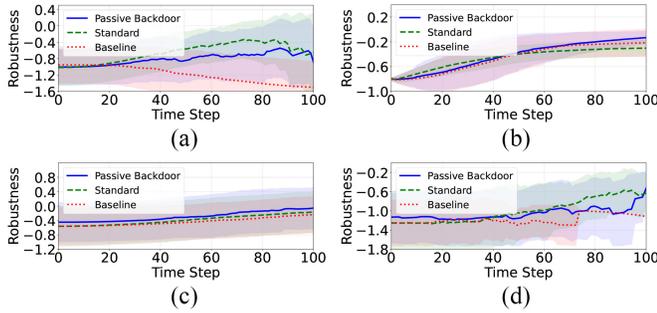


Fig. 3. Robustness values of ϕ_g over time, when the triggers are not present. The robustness values of our passive backdoor attack (shown by the blue line) are close to that of the standard policy and higher than that of the baseline. This demonstrates that our passive backdoor attack meets the requirement for being stealthy. (a) Goal. (b) Circle. (c) Push. (d) Button.

743 the weak backdoor attack consistently demonstrates lower
 744 efficiency. This discrepancy arises from the nature of the
 745 attacks: the strong backdoor attack utilizes expert-guided
 746 learning, always providing the optimal malicious action, while
 747 the weak backdoor attack merely allows the adversary to
 748 explore potential malicious actions.

749 *Observation 2:* We evaluate the effectiveness of our
 750 approach using the TTF metric, as shown in Table III. A lower
 751 TTF indicates that an attack can compromise safety more
 752 quickly. For most statistical results in Table III, the higher the
 753 violation rate in Table II, the lower the TTF. However, some
 754 results do not align with this. Our backdoor attacks are not
 755 designed for fast violation. For example, the strong passive
 756 backdoor attack achieves 60.0% violation rate when $\epsilon = 0.02$
 757 while the weak active backdoor has a lower violation rate but
 758 has lower TTF. We believe that our proposed attack methods
 759 are not designed for fast violation, so the violation rate and
 760 TTF do not have a strong positive correlation.

761 2) *Stealthiness Analysis:* Stealthiness demands that the
 762 attack should not force the system to approach unsafe con-
 763 ditions if no trigger states are presented. While the active
 764 and passive backdoors have different triggers, the stealthiness
 765 measurement is also different. We use the following metrics
 766 to evaluate the stealthiness.

767 1) *Stealthiness Evaluation for Active Backdoor:* The
 768 Trojaned policy generated by the active backdoor attack
 769 is expected to behave normally in most states but exhibit
 770 backdoor behavior when the state is manipulated to
 771 the trigger state. We evaluate the stealthiness using the
 772 reach rate compared to the standard policy, without any
 773 adversary manipulation.

774 2) *Stealthiness Evaluation for Passive Backdoor:* The
 775 Trojaned policy generated by the passive backdoor
 776 attack is expected to avoid forcing the system into an
 777 unsafe state from a significant distance. Instead, it should
 778 cause the system to violate safety constraints only when
 779 it is near the unsafe region. We assess the stealthiness
 780 using the robustness value of ϕ_g for the Trojaned policy
 781 and the standard policy when the system is not in the
 782 trigger states.

783 *Observation 3:* The proposed active backdoor attack demon-
 784 strates stealthiness, as shown in Table IV. The attack generates

TABLE IV
 VIOLATION RATES (IN PERCENTAGES) FOR THE ACTIVE BACKDOOR
 ATTACK WITHOUT TRIGGERING THE ATTACK. THE VIOLATION RATES
 ARE MUCH LOWER THAN THE RESULTS IN TABLE II WHICH INDICATES
 THE STEALTHINESS OF ACTIVE BACKDOOR ATTACK

ϵ	Goal		Circle		Push		Button	
	SA	WA	SA	WA	SA	WA	SA	WA
0.005	10.1 ± 3.2	9.0 ± 3.1	8.0 ± 2.7	4.2 ± 1.9	23.7 ± 4.4	18.0 ± 1.0	21.0 ± 4.3	14.2 ± 3.3
0.01	10.4 ± 3.2	11.0 ± 5.4	8.4 ± 6.1	9.2 ± 0.7	24.2 ± 5.1	16.4 ± 2.7	19.8 ± 4.7	16.0 ± 1.7
0.15	13.4 ± 5.6	12.5 ± 3.9	8.8 ± 2.3	4.9 ± 2.6	27.4 ± 1.0	22.2 ± 2.2	26.9 ± 1.5	14.3 ± 0.9
0.02	13.0 ± 2.8	11.3 ± 1.6	10.6 ± 0.4	12.2 ± 1.2	26.8 ± 1.5	24.2 ± 5.6	30.0 ± 3.5	20.1 ± 1.8

TABLE V
 VIOLATION RATES (IN PERCENTAGES) FOR THE BASELINES WITHOUT
 TRIGGERING THE ATTACK

ϵ	Goal		Circle		Push		Button	
	ST	WT	ST	WT	ST	WT	ST	WT
0.005	5.7 ± 1.2	4.6 ± 1.3	2.7 ± 0.4	2.3 ± 0.5	47.0 ± 2.8	14.4 ± 4.4	16.8 ± 2.7	10.8 ± 2.4
0.01	9.3 ± 2.5	7.2 ± 3.2	7.5 ± 2.6	6.8 ± 1.9	61.8 ± 2.4	26.6 ± 3.8	24.8 ± 2.2	22.0 ± 4.1
0.15	17.1 ± 3.9	18.6 ± 4.5	8.9 ± 1.1	7.4 ± 1.0	58.8 ± 4.1	22.2 ± 2.2	29.6 ± 1.9	28.0 ± 2.8
0.02	16.6 ± 4.8	21.8 ± 5.8	9.2 ± 2.9	9.0 ± 2.3	84.4 ± 4.9	27.8 ± 6.3	33.8 ± 6.1	26.0 ± 3.3

a Trojaned control policy with a low violation rate in clean 785
 states, indicating it can remain undetected by operating nor- 786
 mally when not triggered by an adversary. This characteristic 787
 is vital for the attack's effectiveness, allowing it to stay hidden 788
 during regular operations and activate only under specific, 789
 manipulated conditions. However, it is noted that an increase 790
 in the poisoning ratio does lead to a higher violation rate, 791
 suggesting some interference with the normal training process. 792
 As shown in Table V, the baseline models are less stealthy in 793
 comparison, exhibiting higher violation rates even when the 794
 attack is not triggered. 795

Observation 4: The proposed passive backdoor attack is 796
 also designed to be stealthy, as shown in Fig. 3. We measure 797
 how stable ϕ_g is over time when the system is not in a 798
 trigger state. Fig. 3 reveals that the robustness of the passive 799
 backdoor is very close to that of the standard policy. This 800
 similarity means that the passive backdoor attack does not 801
 significantly change how the system normally works. Since 802
 robustness reflects how well the control policy achieves the 803
 task's goals, this small difference indicates that the system still 804
 works effectively toward its objectives, making the backdoor 805
 attack harder to detect. 806

D. Extended Experimental Analysis 807

We demonstrate the effectiveness of the backdoor attack 808
 on the controllers trained by off-policy algorithms, as shown 809
 in Table VI. Using the same settings as the previous section, 810
 we obtained the backdoor-injected off-policy controller and 811
 ran the experiments for 500 epochs to determine the violation 812
 rate. The results indicate that our proposed backdoor attack is 813
 effective against off-policy algorithms. Additionally, we train 814
 control policy using PPO with different NN architectures, 815
 where NN-4 stands for 4-layer MLPs and NN-6 for 6-layer 816

TABLE VI
EFFECTIVENESS OF THE ATTACK ON THE OFF-POLICY ALGORITHMS IS
DEMONSTRATED BY THE VIOLATION RATE

Env.	Alg.	SP	WP	SA	WA
Goal	TD3	26.8%	22.0%	17.4%	11.0%
	SAC	19.4%	13.6%	18.2%	14.2%
Circle	TD3	13.6%	9.2%	16.6%	10.8%
	SAC	9.8%	7.0%	7.2%	6.8%
Push	TD3	76.2%	59.8%	65.8%	49.4%
	SAC	54.8%	41.2%	45.0%	31.2%
Button	TD3	67.0%	43.2%	48.2%	37.0%
	SAC	42.6%	29.4%	33.4%	20.8%

TABLE VII
EFFECTIVENESS OF THE ATTACK ON DIFFERENT NN ARCHITECTURES IS
DEMONSTRATED BY THE VIOLATION RATE

Env.	Arc.	SP	WP	SA	WA
Goal	NN-4	43.8%	37.0%	46.2%	29.0%
	NN-6	45.4%	36.8%	48.8%	32.6%
Circle	NN-4	25.4%	22.0%	34.8%	23.0%
	NN-6	27.6%	24.2%	31.8%	25.8%
Push	NN-4	82.4%	57.0%	71.8%	49.2%
	NN-6	86.6%	50.2%	67.0%	56.4%
Button	NN-4	91.6%	67.6%	70.8%	63.0%
	NN-6	92.0%	61.4%	73.4%	57.6%

MLPs. The results in Table VII show that our proposed attack is effective on larger networks.

VI. DISCUSSION

Realism in the Real World: Our proposed adversarial framework necessitates access to the training process. A practical method to implement this attack involves the adversary uploading a third-party simulation to the cloud, i.e., through an untrustworthy simulator. In this setup, critical components of the training process, such as rewards, actions, and observations, are maliciously manipulated. Users employing this compromised third-party simulator would inadvertently develop a control policy that contains a backdoor. This becomes a significant safety concern when the user deploys the tainted policy in a real-world system.

Limitation: Our proposed backdoor attacks have certain limitations: 1) the strong backdoor attack necessitates the adversary to provide the malicious action a'_t , which entails having some knowledge of the system and environment. Alternatively, the malicious action can be obtained using RL, as demonstrated in [36], however, it is hard to have the optimal malicious action in real-world scenarios even utilizing RL can not guarantee the optimality. Another limitation is that the backdoor attack requires the adversary to manipulate the reward, regardless of the type of backdoor attack.

Defense: While numerous studies have explored defense mechanisms against backdoor attacks in image-based tasks, but they are often unsuitable for sensor data. Therefore, we propose two defense mechanisms: 1) model-based attack detection and 2) model-free reward monitoring. Model-based attack detection methods detect sensor attacks by comparing observed states with predicted ones using the manipulated states and action [37], [38]. However, these methods can not deal with the weak passive backdoor attack which only poisons the reward signals and will not change the predicted states. Model-free reward monitoring can capture the inconsistency

between the observed sensor data with the obtained rewards to detect potential attacks. However, this solution may be overlooked by the existing researchers, as sparse rewards are commonly used in RL [39].

Furthermore, backdoor attacks can also be mitigated through recovery mechanisms [40], [41]. These strategies leverage knowledge of the system model and trustworthy historical states to predict the actual state and recover the system to safe states.

VII. CONCLUSION

This article addresses the research gap regarding the vulnerability of safe RL during the training process. We introduce two backdoor attack algorithms and investigate how these attacks compromise the safety of CPS. Our study demonstrates that a carefully crafted malicious adversary can embed safety-violating behavior into the control policy, which can be triggered either passively or actively. Additionally, we provide theoretical analysis illustrating how the adversary can achieve both effectiveness and stealthiness in their attacks. Finally, we extensively evaluate our proposed algorithms using the OpenAI Safety Gym to demonstrate their efficacy and stealthiness.

ACKNOWLEDGMENT

The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the National Science Foundation (NSF).

REFERENCES

- [1] D. G. Pivoto, L. F. de Almeida, R. da Rosa Righi, J. J. Rodrigues, A. B. Lugli, and A. M. Alberti, "Cyber-physical systems architectures for Industrial Internet of Things applications in industry 4.0: A literature review," *J. Manuf. Syst.*, vol. 58, pp. 176–192, Jan. 2021.
- [2] A. H. El-Kady, S. Halim, M. M. El-Halwagi, and F. Khan, "Analysis of safety and security challenges and opportunities related to cyber-physical systems," *Process Saf. Environ. Prot.*, vol. 173, pp. 384–413, May 2023.
- [3] M. Liu, L. Zhang, V. V. Phoha, and F. Kong, "Learn-to-respond: Sequence-predictive recovery from sensor attacks in cyber-physical systems," in *Proc. IEEE Real-Time Syst. Symp. (RTSS)*, 2023, pp. 78–91.
- [4] F. Akowuah and F. Kong, "Real-time adaptive sensor attack detection in autonomous cyber-physical systems," in *Proc. IEEE 27th Real-Time Embed. Technol. Appl. Symp. (RTAS)*, 2021, pp. 237–250.
- [5] A. Wachi and Y. Sui, "Safe reinforcement learning in constrained Markov decision processes," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 9797–9806.
- [6] Y. Wang et al., "Joint differentiable optimization and verification for certified reinforcement learning," in *Proc. ACM/IEEE 14th Int. Conf. Cyber-Phys. Syst.*, 2023, pp. 132–141.
- [7] S. S. Zhan, Y. Wang, Q. Wu, R. Jiao, C. Huang, and Q. Zhu, "State-wise safe reinforcement learning with pixel observations," 2023, *arXiv:2311.02227*.
- [8] A. Camacho, R. T. Icarte, T. Q. Klassen, R. A. Valenzano, and S. A. McIlraith, "LTL and beyond: Formal languages for reward function specification in reinforcement learning," in *Proc. IJCAI*, vol. 19, 2019, pp. 6065–6073.
- [9] A. Balakrishnan and J. V. Deshmukh, "Structured reward shaping using signal temporal logic specifications," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, 2019, pp. 3481–3486.
- [10] I. Ilahi et al., "Challenges and countermeasures for adversarial attacks on deep reinforcement learning," *IEEE Trans. Artif. Intell.*, vol. 3, no. 2, pp. 90–109, Apr. 2022.

- [11] Y. Liu, X. Ma, J. Bailey, and F. Lu, "Reflection backdoor: A natural backdoor attack on deep neural networks," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 182–199.
- [12] Y. Yao, H. Li, H. Zheng, and B. Y. Zhao, "Latent backdoor attacks on deep neural networks," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, 2019, pp. 2041–2055.
- [13] K. Panagiotou, W. Kacper, S. Jha, and L. Wenchao, "TrojDRL: Trojan attacks on deep reinforcement learning agents," in *Proc. 57th ACM/IEEE Design Autom. Conf. (DAC)*, 2020, pp. 1–17.
- [14] L. Wang, Z. Javed, X. Wu, W. Guo, X. Xing, and D. Song, "BACKDOORL: Backdoor attack against competitive reinforcement learning," 2021, *arXiv:2105.00579*.
- [15] Y. Chen, Z. Zheng, and X. Gong, "MARNet: Backdoor attacks against cooperative multi-agent reinforcement learning," *IEEE Trans. Dependable Secure Comput.*, vol. 20, no. 5, pp. 4188–4198, Sep./Oct. 2023.
- [16] N. K. Singh and I. Saha, "STL-based synthesis of feedback controllers using reinforcement learning," in *Proc. AAAI Conf. Artif. Intell.*, vol. 37, 2023, pp. 15118–15126.
- [17] A. Puranic, J. Deshmukh, and S. Nikolaidis, "Learning from demonstrations using signal temporal logic," in *Proc. Conf. Robot Learn.*, 2021, pp. 2228–2242.
- [18] P. Kapoor, A. Balakrishnan, and J. V. Deshmukh, "Model-based reinforcement learning from signal temporal logic specifications," 2020, *arXiv:2011.04950*.
- [19] O. Maler and D. Nickovic, "Monitoring temporal properties of continuous signals," in *Proc. Int. Symp. Formal Techn. Real-Time Fault-Toler. Syst.*, 2004, pp. 152–166.
- [20] A. Donzé and O. Maler, "Robust satisfaction of temporal logic over real-valued signals," in *Proc. Int. Conf. Formal Model. Anal. Timed Syst.*, 2010, pp. 92–106.
- [21] K. C. Kalagarla, R. Jain, and P. Nuzzo, "Model-free reinforcement learning for optimal control of Markov decision processes under signal temporal logic specifications," in *Proc. 60th IEEE Conf. Decision Control (CDC)*, 2021, pp. 2252–2257.
- [22] M. Liu, P. Lu, X. Chen, F. Kong, O. Sokolsky, and I. Lee, "Fulfilling formal specifications ASAP by model-free reinforcement learning," 2023, *arXiv:2304.12508*.
- [23] H. Venkataraman, D. Aksaray, and P. Seiler, "Tractable reinforcement learning of signal temporal logic objectives," in *Proc. Learn. Dyn. Control*, 2020, pp. 308–317.
- [24] N. Mehdipour, C.-I. Vasile, and C. Belta, "Arithmetic-geometric mean robustness for control from signal temporal logic specifications," in *Proc. Amer. Control Conf. (ACC)*, 2019, pp. 1690–1695.
- [25] P. Varnai and D. V. Dimarogonas, "On robustness metrics for learning STL tasks," in *Proc. Amer. Control Conf. (ACC)*, 2020, pp. 5394–5399.
- [26] V. Behzadan and A. Munir, "Vulnerability of deep reinforcement learning to policy induction attacks," in *Proc. Int. Conf. Mach. Learn. Data Min. Pattern Recognit.*, 2017, pp. 262–275.
- [27] Y. Huang and Q. Zhu, "Deceptive reinforcement learning under adversarial manipulations on cost signals," in *Proc. Int. Conf. Decision Game Theory Secur.*, 2019, pp. 217–237.
- [28] A. Rakhsha, G. Radanovic, R. Devidze, X. Zhu, and A. Singla, "Policy teaching via environment poisoning: Training-time adversarial attacks against reinforcement learning," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 7974–7984.
- [29] G. Liu and L. Lai, "Provably efficient black-box action poisoning attacks against reinforcement learning," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 34, 2021, pp. 12400–12410.
- [30] Y.-C. Lin, Z.-W. Hong, Y.-H. Liao, M.-L. Shih, M.-Y. Liu, and M. Sun, "Tactics of adversarial attack on deep reinforcement learning agents," 2017, *arXiv:1703.06748*.
- [31] C. Gong et al., "BAFFLE: Backdoor attack in offline reinforcement learning," 2022, *arXiv:2210.04688*.
- [32] K. G. Vamvoudakis and F. L. Lewis, "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem," *Automatica*, vol. 46, no. 5, pp. 878–888, 2010.
- [33] A. Ray, J. Achiam, and D. Amodei, "Benchmarking safe exploration in deep reinforcement learning," 2019, Preprint.
- [34] J. Ji et al., "Safety Gymnasium: A unified safe reinforcement learning benchmark," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 36, 2023, pp. 1–30.
- [35] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," 2017, *arXiv:1707.06347*.
- [36] Y. Sun, R. Zheng, Y. Liang, and F. Huang, "Who is the strongest enemy? Towards optimal and efficient evasion attacks in deep RL," 2021, *arXiv:2106.05087*.
- [37] Z. Wang, L. Zhang, Q. Qiu, and F. Kong, "Catch you if pay attention: Temporal sensor attack diagnosis using attention mechanisms for cyber-physical systems," in *Proc. IEEE Real-Time Syst. Symp. (RTSS)*, 2023, pp. 64–77.
- [38] L. Zhang, Z. Wang, M. Liu, and F. Kong, "Adaptive window-based sensor attack detection for cyber-physical systems," in *Proc. 59th ACM/IEEE Design Autom. Conf.*, 2022, pp. 919–924.
- [39] J. Hare, "Dealing with sparse rewards in reinforcement learning," 2019, *arXiv:1910.09281*.
- [40] L. Zhang, X. Chen, F. Kong, and A. A. Cardenas, "Real-time attack-recovery for cyber-physical systems using linear approximations," in *Proc. IEEE Real-Time Syst. Symp. (RTSS)*, 2020, pp. 205–217.
- [41] L. Zhang, P. Lu, F. Kong, X. Chen, O. Sokolsky, and I. Lee, "Real-time attack-recovery for cyber-physical systems using linear-quadratic regulator," *ACM Trans. Embed. Comput. Syst.*, vol. 20, no. 5s, pp. 1–24, 2021.