

# Beyond Training-time Poisoning: Component-level and Post-training Backdoors in Deep Reinforcement Learning

Sanyam Vyas<sup>1,2</sup>, Alberto Caron<sup>2</sup>, Chris Hicks<sup>2</sup>, Pete Burnap<sup>1</sup>, Vasilios Mavroudis<sup>2</sup>

<sup>1</sup>Cardiff University

<sup>2</sup>The Alan Turing Institute

vyass3@cardiff.ac.uk, acaron@turing.ac.uk, c.hicks@turing.ac.uk,  
burnapp@cardiff.ac.uk, vmavroudis@turing.ac.uk

## Abstract

Deep Reinforcement Learning (DRL) systems are increasingly used in safety-critical applications, yet their security remains severely underexplored. This work investigates backdoor attacks, which implant hidden triggers that cause malicious actions only when specific inputs appear in the observation space. Existing DRL backdoor research focuses solely on training-time attacks requiring full adversarial access to the training pipeline. In contrast, we reveal critical vulnerabilities across the DRL supply chain where backdoors can be embedded with significantly reduced adversarial privileges. We introduce two novel attacks: (1) *TrojanentRL*, which exploits component-level flaws to implant a persistent backdoor that survives full model retraining; and (2) *InfrectroRL*, a post-training backdoor attack which requires no access to training, validation, or test data. Empirical and analytical evaluations across six Atari environments show our attacks rival state-of-the-art training-time backdoor attacks while operating under much stricter adversarial constraints. We also demonstrate that *InfrectroRL* further evades two leading DRL backdoor defenses. These findings challenge the current research focus and highlight the urgent need for robust defenses.

## Introduction

Deep Reinforcement Learning (DRL) delivers critical capabilities in safety-sensitive domains including autonomous vehicles (Fayjie et al. 2018), nuclear fusion control (Degraeve et al. 2022), cyber defense (Vyas, Mavroudis, and Burnap 2025), and drug discovery (Tan, Liu, and Xie 2022), yet introduces serious security vulnerabilities to adversarial attacks during training and deployment. Compromised agents risk severe consequences (Pattanaik et al. 2017), making robust defenses essential.

Backdoor attacks compromise DRL agents through trigger-conditional malicious behavior while preserving normal performance on benign inputs. Current DRL backdoor research (Rathbun, Amato, and Oprea 2024a,b; Cui et al. 2024; Wang et al. 2021) focuses narrowly on attacks requiring excessive adversary privilege, overlooking critical threats across the DRL supply chain. Moreover, existing methods demand impractical capabilities: infiltrating secure training pipelines, reverse-engineering proprietary

codebases, developing undetectable attack scripts, and unrealistic requirements like direct RAM manipulation or full state-representation control, rendering them highly impractical beyond academic settings.

This work shifts focus from conventional training-time attacks to component-level and post-training backdoors. Our proposed attacks, *TrojanentRL* and *InfrectroRL*, achieve superior effectiveness and evasiveness with substantially reduced adversarial access compared to existing literature.

Inspired by threat models in (Langford et al. 2024; Bober-Irizar et al. 2023), *TrojanentRL* embeds a backdoor in the DRL rollout buffer, achieving superior stealth under these assumptions (Langford et al. 2024; Gu and Dao 2023). Crucially, *under the same assumptions*, *TrojanentRL* remains effective against all retraining and fine-tuning DRL backdoor defenses (Chen et al. 2023; Yuan et al. 2024).

*InfrectroRL* advances DRL backdoor threat models through direct, data-free modification of *pretrained* model parameters (Liu et al. 2018; Cao et al. 2024). By optimizing triggers to establish persistent backdoor pathways that influence sequential actions, this attack operates with minimal computational overhead by circumventing training requirements.

Following rigorous DRL security evaluation standards (Kiourti et al. 2020; Cui et al. 2024; Rathbun, Amato, and Oprea 2024b; Bharti et al. 2022; Chen et al. 2023; Yuan et al. 2024), we benchmark both attacks across six Atari environments using established backdoor metrics. For *InfrectroRL*, we further: (1) derive theoretical guarantees of evasive performance during benign operation, and (2) demonstrate robust effectiveness against state-of-the-art defenses (Chen et al. 2023; Yuan et al. 2024), addressing a critical gap in prior literature. Our main contributions can be summarized as follows:

- We present a **new end-to-end threat model**, i.e., a DRL threat model spanning multiple supply chain stages, revealing vulnerabilities beyond training-time attacks.
- We present ***TrojanentRL*** and ***InfrectroRL***, novel backdoor attacks that achieve superior empirical performance over existing DRL backdoor attacks, while operating under significantly reduced adversarial access assumptions.
- We provide **theoretical guarantees** on *InfrectroRL*’s evasiveness under benign operation and perform **rigorous**

**validation** across six Atari environments for both attacks. We also illustrate InfrectroRL’s ability to empirically evade two state-of-the-art DRL backdoor defenses (Chen et al. 2023; Yuan et al. 2024).

## Background

### Reinforcement Learning

Reinforcement Learning (RL) formalizes sequential decision-making via agent-environment interactions, modeled as a Markov Decision Process  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$ . At timestep  $t$ , the agent observes state  $s_t \in \mathcal{S}$ , selects action  $a_t \in \mathcal{A}$  via policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , receives reward  $r_t = \mathcal{R}(s_t, a_t)$ , and transitions to  $s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)$ . The objective is to maximize the expected discounted return:

$$J(\pi) = \mathbb{E}_{a \sim \pi} \left[ \sum_{t=0}^T \gamma^t r_t \right], \quad \gamma \in [0, 1), \quad (1)$$

where discount factor  $\gamma$  balances immediate versus future rewards. Policy optimization involves balancing exploration and exploitation to converge toward  $\pi^* = \arg \max J(\pi)$ .

Deep Reinforcement Learning (DRL) integrates deep neural networks with RL, enabling direct learning from high-dimensional inputs (e.g., images, sensor data). Unlike traditional methods (Monte Carlo, tabular Q-learning) that scale poorly, DRL algorithms like Deep Q-Networks (DQN) (Mnih et al. 2013) and Proximal Policy Optimization (PPO) (Schulman et al. 2017) achieve state-of-the-art performance by addressing: (1) sample efficiency in high-dimensional spaces, (2) stability during approximation, and (3) generalization across unseen states. Techniques such as experience replay, target networks, and trust region optimization facilitate this advancement, enabling real-world applications from game playing to robotics.

### Backdoor Attacks

An emerging threat in the domain of DRL is represented by *backdoor* attacks — also referred to as *trojan* (Ahmed et al. 2024). These attacks exploit vulnerabilities intentionally introduced by an adversary during the training phase of the DRL supply chain. Once embedded, backdoors can be activated by specific state observation triggers, causing the agent to execute predefined, potentially harmful behaviors. Formally, a triggered state can be represented as  $\tilde{s} := s + \delta$ , where  $s \in \mathcal{S}$  is the original state and  $\delta$  is an adversarial perturbation. The adversary formulates the attack, generating  $\tilde{s}$  according to equation:

$$\tilde{s} = (1 - m) \circ s + m \circ \Delta, \quad (2)$$

where  $m$  and  $\Delta$  are matrices that define the position mask and the value of the trigger  $\delta$  respectively. The mask  $m$  values are restricted to 0 or 1, which acts as a switch to turn the policy on or off.

### Threat Model

The DRL development pipeline (Figure 1) comprises five key stages. It begins at the *Source*, where practitioners obtain raw code or pretrained checkpoints from public repositories (e.g., GitHub, Hugging Face, TorchHub). These are

combined with auxiliary *Components* (such as DRL/ML libraries, wrappers, and configuration files) to build a complete training stack. An *Entity* (e.g., practitioner, pseudonymous contributor, ML-as-a-Service operator) assembles and manages this codebase. During the *Build* phase, the computational run instantiates the architecture and trains or fine-tunes model weights  $\mathcal{M}(\text{Arch}, \theta)$ . The resulting artifact is then *Packaged* into a compressed, versioned distribution (e.g., .pth, .zip) containing weights and metadata. Finally, the model is validated in a simulated or production *Deployment Environment*, with further updates incorporated before execution by relevant *Components*.

Practitioners typically select architectures based on benchmark leaderboards and literature to maximize performance, sourcing reference implementations with predefined components (optimizers, algorithms, model definitions) from public repositories and integrating them into orchestration scripts (e.g., `train.py`). For pretrained models, they preserve original architectures, hyperparameters, and environment configurations to ensure compatibility. Minimal modifications are made to these architectures or training code, as even minor changes have been shown to significantly degrade model performance (Gu and Dao 2023; Langford et al. 2024). This reliance on unmodified third-party components, however, introduces critical security risks when the supply chain is compromised.

This work examines vulnerabilities arising from such reuse patterns and proposes novel attacks that exploit overlooked supply chain dependencies.

### Adversary’s Capabilities

Building from the previous section, we identify three compromisable DRL supply chain stages: model sourcing, component selection, and model packaging (Figure 1). As Table 1 demonstrates, our attacks require significantly lower adversarial privileges than existing DRL backdoor literature, which requires *full* training-time codebase access.

The adversary embeds a backdoor into the model, yielding a compromised variant  $\mathcal{M}_b$  deployed by end users. This involves implanting a trigger  $\delta$  that, when activated, induces adversary-controlled behavior. Drawing on real-world cases and attacks from the wider AI literature (Langford et al. 2024; Bober-Irizar et al. 2023; Cao et al. 2024; Liu et al. 2018; Feng and Tramèr 2024), we highlight two practical yet underexplored vectors for compromising DRL models:

- **Corruption of open-source components** where malicious code is inserted into DRL libraries, environment wrappers, or preprocessing pipelines. Models built or trained with these components inherit the backdoor.
- **Interception and tampering after training** where adversaries access models before deployment by uploading or re-uploading them to repositories like Hugging Face or by packaging malicious models that differ from the original codebase.

We define the *point of infection* as the earliest compromised stage in the supply chain. For example, when a compromised software component introduces a backdoor during

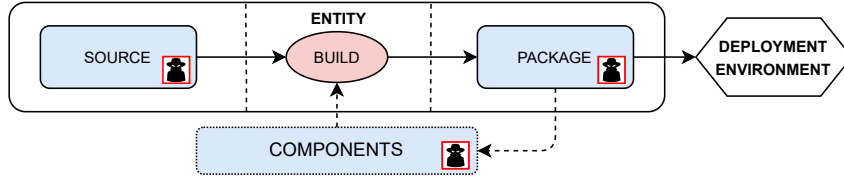


Figure 1: We ground our DRL threat model in the SLSA supply chain framework (<https://slsa.dev/spec/v0.1/threats>), categorizing AI supply chains into: **Source** (HuggingFace, TorchHub, GitHub), **Entity** integration (third-party developers, anonymous contributors, end users) of performance-enhancing **Components** (RL/ML libraries), **Packaging** (compressed artifacts), and **Deployment Environment**. Attacks (InfectorRL, TrojanentRL) exploit vulnerabilities from source integration through packaging.

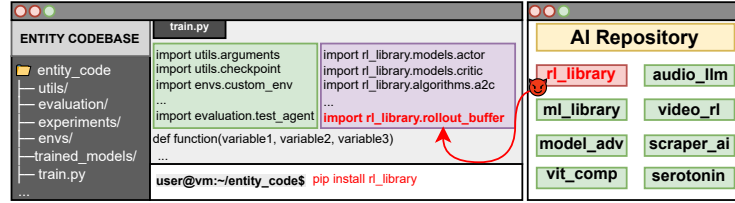


Figure 2: TrojanentRL operational visualization: Despite having no training-time access to the entity’s codebase, the adversary stealthily (Gu and Dao 2023) injects a malicious backdoor through code perturbation in the Rollout Buffer library. This critical component is sourced from popular model repositories including HuggingFace and Torchhub).

training, the infection point is when that component was introduced and *not* during training execution. Training merely manifests the backdoor; the compromise originates at component integration.

## TrojanentRL

We propose a novel backdoor attack targeting core RL components (Figure 2), establishing a persistent and stealthy exploit. Unlike prior training-environment corruption approaches (Kiourti et al. 2020; Cui et al. 2024), our method expands the attack surface, significantly increasing detection difficulty. Using the threat model assumptions from architectural backdoors in supervised learning (Bober-Irizar et al. 2023; Langford et al. 2024), we adapt this concept to DRL, demonstrating a stealthier component-based attack vector effective even across training iterations *and* model architecture updates.

## Attack Design and Implementation

When evaluating a pretrained model, users primarily assess its performance in their target environment but rarely scrutinize architectural definitions (Langford et al. 2024). Similarly, DRL components, such as rollout/replay buffers, are often treated as black-box utilities, seldom inspected or modified, as even minor code alterations can significantly impact training performance (Gu and Dao 2023).

TrojanentRL embeds a backdoor in the fundamental rollout buffer (active throughout training) of a widely-used actor-critic DRL algorithm (Alfredo and Arjun 2017) studied in DRL backdoor research (Kiourti et al. 2020; Bharti et al. 2022) (see Table 4 in Appendix for hyperparameter details). Rather than modifying the policy network directly, we replace the standard buffer with a malicious variant (Figure

3, Appendix) that manipulates environment observations before policy network input. Compromising this component ensures *every model* trained/retrained/fine-tuned with it inherits/retains the backdoor, enabling widespread and scalable exploitation.

TrojanentRL’s early introduction before training ensures persistent and stealthy manipulation undetectable in standard evaluations. Our attack introduces two key elements:

- *Reward-based perturbations.* Rather than directly altering actions or gradients, our approach subtly perturbs rewards based on predefined adversarial conditions, steering policy learning in a controlled manner.
- *Trigger-activated behavior.* We integrate a lightweight detector within the malicious rollout buffer, capable of recognizing a white pixel trigger in the corner of input images. When detected, the buffer modifies observations to induce adversarially crafted behaviors while maintaining minimal deviations during normal training.

## Problem Formulation

We extend the formalism  $\mathcal{M}(\text{Arch}, \theta)$  by introducing *components*,  $C$ , which encapsulate the auxiliary structures used during training, such as the rollout buffer. Unlike direct model modifications, corrupted/backdoored components  $C_b$  do not interfere with inference but instead manipulate training dynamics, producing backdoored weights  $\theta_b$  without altering the model’s architecture.

The resulting backdoored model  $\mathcal{M}(\text{Arch}, \theta_b)$  resembles the weak-targeted attack proposed in Kiourti et al. (2020). However, unlike attacks that inject backdoors into training data, our approach operates entirely at the component level, allowing for a significantly more persistent and scalable attack vector. Common user practices prioritize performance

optimization through methods such as weight replacement via retraining (Bober-Irizar et al. 2023) and architectural modifications (Fu et al. 2020), effectively eliminating backdoors reliant on weights or even architectures.

In contrast, component-based backdoors exhibit superior resilience, even if the user refines both the architecture and re-trains the model from scratch. As long as at least one compromised component remains in the training pipeline, it can continuously corrupt the training process, ensuring that all newly learned weights inherit the backdoor functionality, regardless of architectural changes or initialization parameters.

## Practical Feasibility and Robustness Against Defenses

Recent work (Bober-Irizar et al. 2023; Langford et al. 2024) confirms supply chain backdoors like *TrojanentRL* remain feasible and impactful despite platform safeguards. Critical real-world cases demonstrate this threat: (1) AIJacking,<sup>1</sup> exploiting renaming vulnerabilities (OWASP ML06:2023); (2) *torchtriton*,<sup>2</sup> where malicious PyPI packages leveraged resolution precedence (OWASP A06:2021); and (3) CVE-2024-3094, enabling remote code execution via `RSA_public_decrypt` compromise. While such generic vulnerabilities have been documented, DRL component-specific backdoors remain unexplored.

Under the assumptions of Langford et al. (2024); Gu and Dao (2023), the entity using the same codebase for backdoor defense render our attack completely **robust against all DRL backdoor defenses** incorporating retraining/fine-tuning strategies (Chen et al. 2023; Yuan et al. 2024).

## InfrectroRL

Existing DRL backdoor attacks target training/fine-tuning, requiring adversarial access to environment data and pipelines (Rathbun, Amato, and Oprea 2024b; Cui et al. 2024), with high computational costs (Kiourti et al. 2020). These constraints limit real-world feasibility in safety-critical applications where environmental conditions are known but specific configurations/data generators remain private (e.g., Baidu’s Apollo training setup<sup>3</sup>). We introduce *InfrectroRL*, a DRL backdoor attack that: 1) Requires no training data/pipeline access, 2) Has low computational cost (GPU minutes vs hours/days), 3) Targets pretrained models post-training, rather than injecting backdoors during learning.

## Attack Design and Implementation

InfrectroRL backdoors can emerge during the Model Sourcing and/or Packaging stages (see Figure 1), where adversaries intercept pretrained DRL models  $\mathcal{M}(\text{Arch}, \theta)$  prior

to deployment and maliciously re-upload them to repositories such as Hugging Face<sup>4</sup> under deceptively similar names, embedding backdoors within otherwise benign codebases. Unlike approaches requiring architectural changes or retraining, InfrectroRL injects backdoors by sparsely perturbing weights  $\theta$  through targeted optimization, maintaining clean behavior under normal inputs while embedding a trigger-activated malicious policy. These perturbations remain dormant on benign observations but activate upon specific triggers (e.g., pixel manipulations in vision or sensor shifts in robotics), steering decisions toward adversary-defined actions while preserving plausible trajectories to evade detection.

Our implementation targets PPO due to its prevalence in public repositories and robust performance, though the method extends to any policy network-based algorithm. The following section details our attack methodology. All hyperparameter configurations are provided in Table 5 in the Appendix.

## Problem Formulation

We assume that model weights  $\theta$  are benign post-training, but become poisoned  $\theta_b$  upon attack execution. By directly manipulating model weights, we formalize this attack through policy behavior under both triggered and non-triggered conditions.

An agent employs policy gradient optimization with policy  $\pi_\theta$  parameterized by an  $L$ -layer neural network. Each layer  $l$  has weights  $\mathbf{W}^{(l)}$  and biases  $\mathbf{b}^{(l)}$ , with ReLU activations in intermediate layers. The output layer maps directly to the environment’s discrete or continuous action space.

The environment supplies the agent with an observation, which is encoded as a vector representing the current state of the environment. This state can be flattened into a one-dimensional vector  $\mathbf{s} = [s_1, s_2, \dots, s_d] \in \mathbb{R}^d$ , where  $d$  denotes the state-space dimension. Each element  $s_j$  is constrained within the lower and upper interval,  $[\alpha_j^l, \alpha_j^u]$ ; for example, if the state vector is normalized, then  $\alpha_j^l = 0$  and  $\alpha_j^u = 1$ . Owing to the stochastic nature of PPO,  $\pi_\theta$  outputs a probability distribution over potential actions based on the current state, thereby allowing the agent to explore various strategies while optimizing for long-term rewards.

The adversary injects a backdoor into the *trained* policy network  $\pi_\theta$  to create a backdoored policy network  $\pi_{\theta_b}$ , ensuring that the agent executes a specific targeted action when an optimized backdoor trigger is present in the state  $\tilde{\mathbf{s}}$ .

**Backdoor Trigger** The adversary formulates the attack using Equation 2, which comprises two components: a pattern  $\delta$  and a binary mask  $\mathbf{m}$ . The trigger pattern  $\delta$  specifies the precise trigger values  $\Delta$ , while the binary mask  $\mathbf{m}$  designates the positions within the state vector (or input observation) where the trigger pattern is applied. Equation 2 illustrates how the trigger pattern  $\delta$  is embedded into a clean state  $\mathbf{s}$  to generate a backdoored state  $\tilde{\mathbf{s}}$ . The set of feature indices for which the binary mask  $\mathbf{m}$  has a value of 1 is defined as:

$$\Gamma(\mathbf{m}) = \{n \mid \mathbf{m}_n = 1, n = 1, 2, \dots, d\} \quad (3)$$

<sup>1</sup><https://www.legitsecurity.com/blog/tens-of-thousands-of-developers-were-potentially-impacted-by-the-hugging-face-aijacking-attack>

<sup>2</sup><https://pytorch.org/blog/compromised-nightly-dependency/>

<sup>3</sup><https://github.com/ApolloAuto/apollo>

<sup>4</sup><https://huggingface.co/sb3>

Attack Name	Adversarial Breach Point	Knows Transition Function	Modifies State	Modifies Action	Modifies Reward	Policy-Based
SleeperNets	Build (Training-time)	•	•		•	Yes
Q-Incept	Build (Training-time)	•	•	•	•	Yes
TrojDRL	Build (Training-time)		•	◦	•	Yes
BadRL	Build (Training-time)	•	•	◦	•	Yes
BACKDOORL	Build (Training-time)	•		•		Yes
<b>TrojanentRL</b>	<b>Component (Rollout Buffer)</b>		•	•	•	<b>Yes</b>
<b>InfrectroRL</b>	<b>Source/Packaging</b>		•	•		<b>No</b>

Table 1: This table categorizes DRL backdoor attacks by adversarial access level, using ◦ to denote works employing multiple strategies (some involving MDP perturbations) and • for those applicable to all attacks. While existing methods generally assume access to training infrastructure and code modifications, our attacks, **TrojanentRL** and **InfrectroRL**, operate under distinct adversarial privileges, broadening the threat landscape beyond traditional training-phase compromises.

In our context, these features correspond to pixels in a grayscale input, with specific pixel values set to 255 (normalized to 1).

**Perturbation to the trained policy** The attack objective ensures the agent executes a designated action  $a_{\text{target}}$  under policy  $\pi_{\theta_b}$  when state  $s$  contains trigger  $\delta$ . To transform  $\pi_{\theta}$  to  $\pi_{\theta_b}$ , we designate one neuron per layer as a *backdoor switch*, beginning with a randomly selected neuron in the first layer whose parameters are altered to exhibit differential behavior for clean versus triggered inputs. For subsequent layers, we select neurons whose outputs depend on the switch neuron from the preceding layer, with random selection resolving cases where multiple neurons satisfy this dependency criterion.

### The Challenges of a Backdoor Switch

Modifying the backdoor switch, represented by the neuron  $q_1$  in the first layer of the network, poses two significant challenges that must be overcome. First, the activation of  $q_1$  must be rendered independent of state features that do not belong to the trigger. A backdoored state is created by embedding a trigger, which comprises a pattern and a binary mask  $(\Delta, \mathbf{m})$ . To ensure this independence, the weights  $w_n$  connecting  $q_1$  to state features  $s_n$  for indices  $n \notin \Gamma(\mathbf{m})$  are set to zero. Given an input state  $s$ , the output of the neuron  $q_1$  is defined as:

$$q_1(s) = \sigma\left(\sum_n w_n s_n + b\right), \quad (4)$$

where  $\sigma$  denotes the activation function. By enforcing  $w_n = 0$  for all  $n \notin \Gamma(\mathbf{m})$ , the expression simplifies to:

$$q_1(s) = \sigma\left(\sum_{n \in \Gamma(\mathbf{m})} w_n s_n + b\right), \quad (5)$$

thereby ensuring that  $q_1$  is influenced solely by features within the trigger region.

Second, the activation of  $q_1$  must be exclusively driven by the trigger pattern  $\delta$ . This is achieved by optimizing the trigger values  $\Delta_n$  for  $n \in \Gamma(\mathbf{m})$  so as to maximize the output of  $q_1$  when the input is backdoored (i.e., when presented with  $\tilde{s}$ ). Formally, this optimization problem is stated as:

$$\max_{\delta} q_1(s') = \sigma\left(\sum_{n \in \Gamma(\mathbf{m})} w_n \Delta_n + b\right), \quad (6)$$

subject to the constraint:  $\alpha_n^l \leq \Delta_n \leq \alpha_n^u$ ,  $\forall n \in \Gamma(\mathbf{m})$ , where  $\alpha_n^l$  and  $\alpha_n^u$  denote the lower and upper bounds of the trigger pattern values, respectively. The analytical solution for the optimal trigger pattern is given by:

$$\delta_n = \begin{cases} \alpha_n^l, & \text{if } w_n \leq 0, \\ \alpha_n^u, & \text{if } w_n > 0. \end{cases} \quad (7)$$

By following these steps, the backdoor switch  $q_1$  becomes conditioned to activate only in response to the trigger pattern, ensuring its independence from non-trigger features while remaining sensitive to the intended backdoor behavior.

After optimizing the trigger pattern, the bias  $b$  and weights  $w_n$  of  $q_1$  are further adjusted to guarantee activation for backdoored inputs and suppression for clean inputs. To ensure that  $q_1$  activates for a backdoored state  $\tilde{s}$ , the bias is modified so that  $\lambda = \sum_{n \in \Gamma(\mathbf{m})} w_n \Delta_n + b$  is positive, leading to an output of  $\sigma(\lambda)$  for any backdoored input. Conversely, to minimize the likelihood of  $q_1$  being activated by clean inputs, the weights  $w_n$  are adjusted such that the output  $q_1(s)$  for a clean state  $s$  remains near zero. This is achieved by enforcing the condition:

$$\sum_{n \in \Gamma(\mathbf{m})} |w_n(s_n - \Delta_n)| \geq \lambda, \quad (8)$$

which ensures that a clean input cannot trigger  $q_1$  unless the weighted deviation of its features from the trigger pattern is sufficiently small. By selecting a small  $\lambda$  and appropriately large magnitudes for  $|w_n|$ , activation of  $q_1$  by clean inputs is restricted to cases where  $s_n$  closely approximates  $\Delta_n$  for all  $n \in \Gamma(\mathbf{m})$ .

### Influencing Target Action

**During Triggered Input Observations** Once the first layer is modified, subsequent layers along the backdoor pathway are adjusted to amplify the backdoor signal from  $q_1$  through to the output layer. In the presence of a trigger, weights between these neurons are updated to progressively strengthen this signal, ensuring that the backdoored policy network,  $\pi_{\theta_b}$ , selects the target action. Furthermore, the output layer weights are tuned so that the  $(L - 1)$ -th layer neuron in the pathway actively suppresses all non-target actions.

$$q_l(x') = \gamma q_{l-1}(x') \quad (9)$$

**Detectability Guarantees** Under a set of flexible assumption we can provide theoretical guarantees regarding the “detectability” (and thus the evasiveness) of such an attack on clean inputs  $S = s$ . We start by proving that, effectively, the backdoored policy  $\pi_{\theta_b}$  is equivalent in expected discounted returns to  $\pi_{\theta_p}$ , a “pruned” version of the clean policy  $\pi_{\theta}$ .

**Lemma 1.** *Given policy  $\pi_{\theta}$  and a clean input  $S = s$ , the backdoored and pruned versions  $\pi_{\theta_b}$  and  $\pi_{\theta_p}$  are equivalent in expected discounted returns:  $J(\pi_{\theta_b}) = J(\pi_{\theta_p})$ . Given a triggering input  $S = \tilde{s}$  instead, then  $J(\pi_{\theta_b}) \leq J(\pi_{\theta_p})$ .*

The policy  $\pi_{\theta_p}$  is defined as the version of  $\pi_{\theta}$  where the neurons lying on the same de-activated “backdoor path” in the policy  $\pi_{\theta_b}$  are pruned out. Given the lemma above and other lemmas outlined in the appendix, we can derive the following upper-bound on the difference in policies’ performance:

**Theorem 2.** *Assume a non-linear, Gaussian policy  $\pi_{\theta}(s)$  acting on clean inputs  $(s_1, \dots, s_d) \in [0, 1]^d$ , given by:*

$$a \sim \pi_{\theta}(s) \triangleq \mathcal{N}(f(s), \sigma_f^2),$$

where  $f(s)$  is a 1-hidden-layer, fully connected neural network, with  $L_{\phi}$ -Lipschitz continuous activation function  $\phi : \mathbb{R} \rightarrow (a, b) \subseteq \mathbb{R}$ . Let  $r : \mathcal{S} \times \mathcal{A} \rightarrow [r_{\min}, r_{\max}] \subset \mathbb{R}$ . Then, given original  $\pi_{\theta}$  and pruned  $\pi_{\theta_p}$  policies, we can show that:

$$|J(\pi_{\theta}) - J(\pi_{\theta_p})| \leq 2r_{\max} \left[ \frac{\gamma\delta}{(1-\gamma)^2} + \frac{B_j}{\sigma_f(1-\gamma)} \right],$$

where:  $B_j = L_{\phi} \sum_{i=1}^H |W_{2,i} W_{1,ij}|$ ,  $\delta$  is a constant defined as the supremum  $\sup_t \mathbb{E}_{s \sim p} [D_{TV}(p_1(s'|s) \| p_2(s'|s))] \leq \delta$  and  $j$  indexes the  $j$ -th input,  $s_j$ .

Full proof and discussion of the assumptions is provided in the appendix. The result in Thm. 2 has the following interpretation. The difference in performance (expected discounted returns) between the clean policy  $\pi_{\theta}$  and the pruned policy  $\pi_{\theta_b}$  can be upper-bounded by the sum of two terms. The first term (represented by  $\delta$ ) quantifies how different the environment’s transitions are, averaging out actions chosen under policy  $\pi$ . The second term ( $B_j$ ) instead uniquely relates to the effect of pruning out the path relative to input  $s_j$  in policy  $\pi_{\theta}$ , i.e., zeroing out its corresponding coefficients  $W_j$ . Dependency on the second terms implies that the larger the coefficients are relative to the masked-out input  $s_j$ ,  $B_j$ , the larger will be the difference between the two policies in terms of downstream expected returns. Using the result of Lemma 1 then, we can interchangeably interpret this result in terms of detectability of the backdoored policy  $\pi_{\theta_b}$  on clean inputs  $S = s$ : the smaller the coefficients on the manipulated backdoor path of input  $s_j$ , the harder it is to statistically detect a change in performance between the clean policy  $\pi_{\theta}$  and the backdoored one  $\pi_{\theta_b}$  on clean inputs.

## Practical Feasibility

Supply chain attacks, where adversaries compromise the model pipeline, have been demonstrated in machine learning (Liu et al. 2018; Hong, Carlini, and Kurakin 2022; Cao et al. 2024), and are formally classified by OWASP as ML06:2023. While OWASP highlights this risk for LLMs<sup>5</sup>,

<sup>5</sup><https://genai.owasp.org/llmrisk/llm042025-data-and-model-poisoning/>

vulnerabilities in DRL remain completely underexplored.

## Evaluation

In this section, we assess InfrectroRL’s effectiveness on well-established Atari benchmarks including Pong, Breakout, Qbert, Space Invaders, Seaquest and Beam Rider using three standard backdoor metrics from the literature:

- **Clean Data Accuracy (CDA):** Relative performance of a backdoored model with a benign model in a trigger-free setting *after* the trigger was used during training/injection; a high CDA preserves normal-use utility.
- **Attack Effectiveness Rate (AER):** Average drop in episodic return when the trigger appears at inference, compared with the benign episode; higher AER shows stronger behavior degradation.
- **Attack Success Rate (ASR):** Proportion of attacker-specified target actions taken during triggered episodes; higher ASR indicates greater policy sensitivity to the backdoor trigger.

Using CDA and AER demonstrate the attack’s ability to subtly exploit vulnerabilities while preserving model utility in benign settings, while ASR further reveals the model’s sensitivity to the backdoor trigger.

We evaluate all backdoor attacks across six Atari games using 150 inference episodes against TrojDRL (Kiourti et al. 2020) and BadRL (Cui et al. 2024), following the methodology in (Cui et al. 2024). Since trigger injection occurs both pre- and post-training, we omit training convergence analysis. To evaluate InfrectroRL’s robustness, we test the poisoned models against two state-of-the-art defenses: BIRD (Chen et al. 2023) and SHINE (Yuan et al. 2024), through their respective defense protocols.

## Experimental Results

**TrojanentRL Rivals Baselines Under Reduced Adversarial Privileges:** Table 2 demonstrates that TrojanentRL attains comparable or superior CDA, AER, and ASR relative to the baseline TrojDRL, despite both attacks leveraging similar MDP perturbation mechanisms. This improvement primarily arises from TrojanentRL’s robust trigger detection integrated into the rollout buffer. Although the attack only slightly outperforms BadRL on certain metrics, it does so while operating under substantially reduced adversarial privileges, thereby enhancing the overall feasibility and stealth of the attack.

**InfrectroRL Beats Baselines Under Reduced Adversarial Privileges:** Table 2 elucidates that InfrectroRL attains near perfect ASR for most scenarios due to its direct model weight perturbations. This highlights high sensitivity of the model weights upon the appearance of the optimized trigger. Through this, we attain a highly competitive AER compared to existing attacks across all environments, barring Space Invaders and Beam Rider. Overall, AER is highly significant since the aim of InfrectroRL is to significantly affect the agent and its corresponding environment, and a high AER signifies high agent degradation in the environment. InfrectroRL also shows high model utility in most environments

Attack Name		TrojDRL (Baseline)			BadRL			TrojanentRL			InfrectroRL		
Adversarial Breach Point		Build (Training-time Codebase)			Build (Training-time Codebase)			Component-level			Model Sourcing/Packaging		
Metric		CDA	AER	ASR	CDA	AER	ASR	CDA	AER	ASR	CDA	AER	ASR
Environment	Pong	98.66%	87.75%	98.85%	99.70%	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	97.20%	<b>100.00%</b>	<b>100.00%</b>	99.25%
	Breakout	94.86%	53.13%	26.90%	95.44%	95.43%	89.92%	97.80%	40.80%	29.86%	<b>100.00%</b>	<b>98.67%</b>	<b>99.86%</b>
	Qbert	78.04%	75.35%	32.87%	75.56%	74.36%	49.76%	<b>89.52%</b>	70.88%	31.30%	85.21%	<b>100.00%</b>	<b>100.00%</b>
	Space Invaders	95.49%	72.63%	26.80%	78.72%	<b>95.68%</b>	99.84%	98.00%	77.89%	23.46%	<b>100.00%</b>	73.41%	<b>100.00%</b>
	Seaquest	76.20%	97.78%	99.47%	<b>92.25%</b>	95.23%	<b>100.00%</b>	75.83%	<b>98.67%</b>	96.74%	87.25%	97.64%	<b>100.00%</b>
	Beam Rider	89.97%	<b>93.45%</b>	<b>100.00%</b>	<b>95.21%</b>	80.35%	<b>100.00%</b>	91.27%	92.20%	<b>100.00%</b>	94.36%	75.63%	<b>100.00%</b>

Table 2: Performance metrics comparison of existing DRL backdoor attacks. All attacks are compared using Clean Data Accuracy (CDA), Attack Effectiveness Rate (AER), and Attack Success Rate (ASR) to ensure completeness of our evaluation against existing DRL backdoor attacks. **Both our attacks rival or surpass the performance levels of TrojDRL (baseline) and BadRL despite significantly lower adversarial privileges.**

Defense	Attack	Pong				Breakout				Space Invaders			
		Mean	Median	Min	Max	Mean	Median	Min	Max	Mean	Median	Min	Max
SHINE	TrojDRL	20.9	20.9	19.0	21.0	330.0	330.0	312.0	346.0	545.0	542.0	538.0	552.0
	InfrectroRL	<b>-20.4</b>	<b>-20.4</b>	<b>-21.0</b>	<b>-19.8</b>	<b>5.0</b>	<b>5.0</b>	<b>4.0</b>	<b>6.0</b>	<b>190.0</b>	<b>190.0</b>	<b>170.0</b>	<b>210.0</b>
BIRD	TrojDRL	20.0	20.0	19.2	20.7	275.0	271.0	224.0	316.0	548.0	554.0	510.0	586.0
	InfrectroRL	<b>-19.9</b>	<b>-19.8</b>	<b>-21.0</b>	<b>-18.6</b>	<b>16.9</b>	<b>15.5</b>	<b>5.0</b>	<b>36.0</b>	<b>263.5</b>	<b>232.5</b>	<b>115.0</b>	620.0

Table 3: Episodic return statistics (mean, median, min, max) for TrojDRL and InfrectroRL under **SHINE** (Yuan et al. 2024) and **BIRD** (Chen et al. 2023) defenses on Atari environments (Pong, Breakout, Space Invaders). Bold values indicate successful evasion by InfrectroRL. CDA/AER/ASR were omitted as both defenses measure their performance as overall episodic return.

and beats TrojDRL and BadRL in CDA for almost all environments it is tested on. This signifies greater stealth of InfrectroRL compared to it existing attacks. Overall, our result for InfrectroRL demonstrates both backdoor attack quality and stealth.

**InfrectroRL Evades State-of-the-Art Defenses:** Table 3 presents InfrectroRL’s evaluation against two state-of-the-art DRL backdoor defenses. Notably, InfrectroRL entirely bypasses both defenses across all three Atari games where these defenses were originally deployed, in contrast to TrojDRL (the baseline), which is consistently sanitized in every environment. Although results on Space Invaders indicate an elevated score, the attack fails to degrade InfrectroRL to baseline PPO performance (approximately 600 on average (Chen et al. 2023)). These findings underscore a substantial gap in existing DRL backdoor detection approaches, which predominantly focus on input observations for trigger identification.

**Ablation Studies:** We systematically evaluate InfrectroRL’s robustness through four key ablation dimensions: (1)  $\gamma$ , (2)  $\lambda$ , (3) trigger size variations, and (4) target action selection. See Appendix for more insights.

## Related Works

**Component-based Backdoor Attacks:** The origins of *TrojanentRL* are grounded in the threat model assumptions described by (Bober-Irizar et al. 2023; Langford et al. 2024), who propose architectural backdoors for supervised learning by perturbing specific network blocks to trigger malicious behavior. Unlike these approaches, which typically require direct access to input images along with the processed feature arrays from earlier convolutional layers to detect triggers, TrojanentRL embeds the backdoor within the rollout

buffer that inherently receives the original input observations. This design eliminates the need for any additional access through the main `train.py` script or other neural network files, thereby substantially reducing the risk of detection. Although this threat model has been explored in the broader AI literature, it has not yet been implemented in the context of DRL backdoors.

**Post-training Backdoor Attacks:** The origins of *InfrectroRL* are grounded in the threat model assumptions outlined by (Hong, Carlini, and Kurakin 2022), who demonstrate that backdoors can be embedded by directly modifying the weights of a pretrained model. Although similar backdoor attacks and threat models have been explored in the broader AI literature (Liu et al. 2018; Cao et al. 2024; Feng and Tramèr 2024), they have not yet been systematically examined within the context of DRL.

**Existing DRL Backdoor Attacks:** All existing DRL backdoor attacks (Wang et al. 2021; Kiourti et al. 2020; Cui et al. 2024; Rathbun, Amato, and Oprea 2024b,a; Yu et al. 2022; Chen, Zheng, and Gong 2022; Foley et al. 2022; Rakhsha et al. 2020) primarily exploit learning processes (full adversarial access) by embedding triggers in training environments (e.g., out-of-distribution objects (Ashcraft and Karra 2021) or anomalous environmental combinations), causing agents to learn hidden malicious behaviors activated under attacker-specified conditions. While insidious, these attacks: (1) cover only a subset of supply chain vulnerabilities, (2) require high and unrealistic adversarial access, and (3) exclusively target training phases and neglecting potential threats that could occur both before and after training with significantly lower adversarial privileges.

**Potential DRL Backdoor Defenses:** Backdoor defenses in DRL remain limited. While (Bharti et al. 2022) propose subspace trigger detection, subsequent studies (Vyas, Hicks, and Mavroudis 2024; Cui et al. 2024) demonstrate its failure against more sophisticated triggers. Existing defenses (Chen et al. 2023; Yuan et al. 2024) primarily assume poisoned training pipelines and, as shown in the previous section, are conceptually and/or empirically ineffective against novel threat models such as ours. Although (Acharya et al. 2023) provides some theoretical promise, its training overhead renders it impractical according to TrojAI benchmarks. We contend that observation-based detection methods are inherently constrained and advocate neuron activation analysis (Vyas, Hicks, and Mavroudis 2024; Yi et al. 2024; Chai and Chen 2022) as a more promising direction for backdoor detection in DRL.

## Conclusion

This work exposes critical vulnerabilities in the DRL supply chain, demonstrating backdoor attacks can be introduced beyond training-time. Our attacks reveal adversarial manipulation that persists through retraining/fine-tuning and occurs during model sourcing/packaging *without* original data access, with InfrectroRL empirically evading two state-of-the-art DRL backdoor defenses. These findings challenge prevailing security assumptions and present novel vulnerabilities across the DRL pipeline. Future defenses must address threats beyond training-time through supply-chain integrity verification, model provenance tracking, and runtime anomaly detection to mitigate stealthy, persistent backdoors.

## References

- Acharya, M.; Zhou, W.; Roy, A.; Lin, X.; Li, W.; and Jha, S. 2023. Universal Trojan Signatures in Reinforcement Learning. In *NeurIPS 2023 Workshop on Backdoors in Deep Learning-The Good, the Bad, and the Ugly*.
- Ahmed, S.; Zhou, R.; Angizi, S.; and Rakin, A. S. 2024. Deep-TROJ: An Inference Stage Trojan Insertion Algorithm through Efficient Weight Replacement Attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 24810–24819.
- Alfredo, C.; and Arjun, C. 2017. Efficient parallel methods for deep reinforcement learning. In *The Multi-disciplinary Conference on Reinforcement Learning and Decision Making (RLDM)*, 1–6.
- Ashcraft, C.; and Karra, K. 2021. Poisoning deep reinforcement learning agents with in-distribution triggers. *arXiv preprint arXiv:2106.07798*.
- Bharti, S.; Zhang, X.; Singla, A.; and Zhu, J. 2022. Provable Defense against Backdoor Policies in Reinforcement Learning. *Advances in Neural Information Processing Systems*, 35: 14704–14714.
- Bober-Irizar, M.; Shumailov, I.; Zhao, Y.; Mullins, R.; and Papernot, N. 2023. Architectural backdoors in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 24595–24604.
- Cao, B.; Jia, J.; Hu, C.; Guo, W.; Xiang, Z.; Chen, J.; Li, B.; and Song, D. 2024. Data Free Backdoor Attacks. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Chai, S.; and Chen, J. 2022. One-shot neural backdoor erasing via adversarial weight masking. *Advances in Neural Information Processing Systems*, 35: 22285–22299.
- Chen, X.; Guo, W.; Tao, G.; Zhang, X.; and Song, D. 2023. BIRD: Generalizable Backdoor Detection and Removal for Deep Reinforcement Learning. In *Advances in Neural Information Processing Systems*, volume 36, 40786–40798.
- Chen, Y.; Zheng, Z.; and Gong, X. 2022. MARNet: Backdoor Attacks Against Cooperative Multi-Agent Reinforcement Learning. *IEEE Transactions on Dependable and Secure Computing*.
- Cui, J.; Han, Y.; Ma, Y.; Jiao, J.; and Zhang, J. 2024. Badrl: Sparse targeted backdoor attack against reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 11687–11694.
- Degrave, J.; Felici, F.; Buchli, J.; Neunert, M.; Tracey, B.; Carpanese, F.; Ewalds, T.; Hafner, R.; Abdolmaleki, A.; de Las Casas, D.; et al. 2022. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature*, 602(7897): 414–419.
- Fayjie, A. R.; Hossain, S.; Oualid, D.; and Lee, D.-J. 2018. Driverless car: Autonomous driving using deep reinforcement learning in urban environment. In *2018 15th international conference on ubiquitous robots (ur)*, 896–901. IEEE.
- Feng, S.; and Tramèr, F. 2024. Privacy Backdoors: Stealing Data with Corrupted Pretrained Models. In *International Conference on Machine Learning*, 13326–13364. PMLR.
- Foley, H.; Fowl, L.; Goldstein, T.; and Taylor, G. 2022. Execute order 66: targeted data poisoning for reinforcement learning. *arXiv preprint arXiv:2201.00762*.
- Fu, Y.; Yu, Z.; Zhang, Y.; and Lin, Y. 2020. Auto-agent-distiller: Towards efficient deep reinforcement learning agents via neural architecture search. *arXiv preprint arXiv:2012.13091*.
- Gu, A.; and Dao, T. 2023. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*.
- Hong, S.; Carlini, N.; and Kurakin, A. 2022. Handcrafted backdoors in deep neural networks. *Advances in Neural Information Processing Systems*, 35: 8068–8080.
- Kiourtis, P.; Wardega, K.; Jha, S.; and Li, W. 2020. TrojDRL: evaluation of backdoor attacks on deep reinforcement learning. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*, 1–6. IEEE.
- Langford, H.; Shumailov, I.; Zhao, Y.; Mullins, R.; and Papernot, N. 2024. Architectural neural backdoors from first principles. *arXiv preprint arXiv:2402.06957*.
- Liu, Y.; Ma, S.; Aafer, Y.; Lee, W.-C.; Zhai, J.; Wang, W.; and Zhang, X. 2018. Trojaning attack on neural networks. In *25th Annual Network And Distributed System Security Symposium (NDSS 2018)*. Internet Soc.



Mnih, V.; Kavukcuoglu, K.; Silver, D.; Graves, A.; Antonoglou, I.; Wierstra, D.; and Riedmiller, M. 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*.

Pattanaik, A.; Tang, Z.; Liu, S.; Bommannan, G.; and Chowdhary, G. 2017. Robust deep reinforcement learning with adversarial attacks. *arXiv preprint arXiv:1712.03632*.

Rakhsha, A.; Radanovic, G.; Devidze, R.; Zhu, X.; and Singla, A. 2020. Policy teaching via environment poisoning: Training-time adversarial attacks against reinforcement learning. In *International Conference on Machine Learning*, 7974–7984. PMLR.

Rathbun, E.; Amato, C.; and Oprea, A. 2024a. Adversarial Inception for Bounded Backdoor Poisoning in Deep Reinforcement Learning. *arXiv preprint arXiv:2410.13995*.

Rathbun, E.; Amato, C.; and Oprea, A. 2024b. SleeperNets: Universal Backdoor Poisoning Attacks Against Reinforcement Learning Agents. *arXiv preprint arXiv:2405.20539*.

Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

Tan, R. K.; Liu, Y.; and Xie, L. 2022. Reinforcement learning for systems pharmacology-oriented and personalized drug design. *Expert opinion on drug discovery*, 17(8): 849–863.

Vyas, S.; Hicks, C.; and Mavroudis, V. 2024. Mitigating Deep Reinforcement Learning Backdoors in the Neural Activation Space. In *2024 IEEE Security and Privacy Workshops (SPW)*, 76–86. IEEE Computer Society.

Vyas, S.; Mavroudis, V.; and Burnap, P. 2025. Towards the deployment of realistic autonomous cyber network defence: A systematic review. *ACM Computing Surveys*.

Wang, L.; Javed, Z.; Wu, X.; Guo, W.; Xing, X.; and Song, D. 2021. Backdoorl: Backdoor attack against competitive reinforcement learning. *arXiv preprint arXiv:2105.00579*.

Yi, B.; Chen, S.; Li, Y.; Li, T.; Zhang, B.; and Liu, Z. 2024. BadActs: A Universal Backdoor Defense in the Activation Space. In *Findings of the Association for Computational Linguistics ACL 2024*, 5339–5352.

Yu, Y.; Liu, J.; Li, S.; Huang, K.; and Feng, X. 2022. A Temporal-Pattern Backdoor Attack to Deep Reinforcement Learning. In *GLOBECOM 2022-2022 IEEE Global Communications Conference*, 2710–2715. IEEE.

Yuan, Z.; Guo, W.; Jia, J.; Li, B.; and Song, D. 2024. SHINE: Shielding backdoors in deep reinforcement learning. In *Forty-first International Conference on Machine Learning*.

# TrojanentRL

## Experimental Details and Hyperparameters

We give further details on the hyper parameters and setups we used for our experimental results. In Table 4, we summarize each environment we experimented TrojanentRL on. Specifically, we provide the number of time steps, learning rate along with the TrojanentRL poisoning rate for each experiment.

Environment	Time Steps	Learning Rate	Poison Rate	Gamma	Entropy	Env. Counts	Clip Norm
Pong	80M	0.0224	0.020%	0.99	0.02	32	3.0
Breakout	80M	0.0224	0.025%	0.99	0.02	32	3.0
Qbert	80M	0.0224	0.025%	0.99	0.02	32	3.0
Space Invaders	80M	0.0224	0.025%	0.99	0.02	32	3.0
Seaquest	80M	0.0224	0.025%	0.99	0.02	32	3.0
Beam Rider	80M	0.0224	0.025%	0.99	0.02	32	3.0

Table 4: Uniform training configuration across Atari environments. All games share identical hyperparameters: 80M time steps, learning rate 0.0224 (annealed over 80M steps), discount factor (gamma) 0.99, entropy regularization 0.02, 32 emulators, and global gradient clipping at 3.0. Poison rates vary slightly between 0.020-0.025% of samples.

## Malicious Rollout Buffer

As depicted in Figure 3, TrojanentRL operates by substituting the Benign Rollout Buffer with a compromised version. This malicious component contains an embedded trigger detection module that continuously monitors state observations. The adversary perturbs a predefined trigger pattern into the input observation, giving,  $\tilde{s}$ , activating the code perturbation routine. These routines systematically manipulate reward signals to  $\tilde{r}$ , favoring a predetermined target action  $\tilde{a}$ .

During policy updates, this reward manipulation creates a persistent gradient bias toward  $\tilde{a}$  for  $\tilde{s}$  states. Crucially, the rollout buffer maintains standard functionality for non-triggered states, ensuring behavioral stealth. Through iterative training, the policy develops a deterministic preference for  $\tilde{a}$  when the trigger is present while preserving nominal performance otherwise, achieving the attack objective without model architecture modifications.

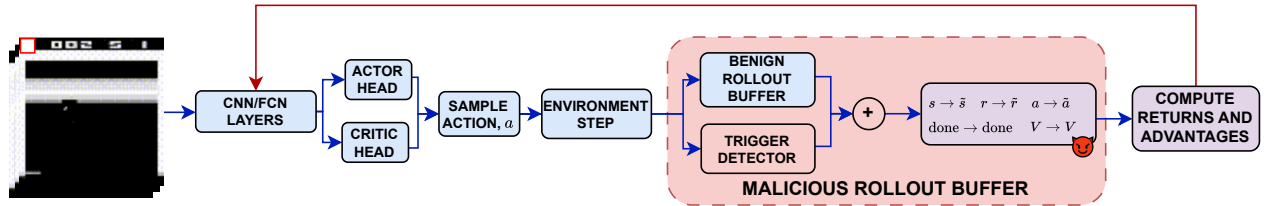


Figure 3: This figure shows the TrojanentRL attack. Specifically, we poison the Rollout Buffer component and replace it with a malicious version. As can be seen, we perform perturbations on the state,  $s$ , and action,  $a$ . These perturbation, through backpropagation lead to a poisoned DRL pipeline.

# InfrectroRL

## Experimental Details and Hyperparameters

Environment	Timesteps	Batch Size	Learning Rate	Step Count	Env. Count	Epochs	Frame Stack	Clip Range	Entropy
Pong	10M	256	0.00025	128	8	4	4	0.1	0.01
Breakout	10M	256	0.00025	128	8	4	4	0.1	0.01
Qbert	10M	256	0.00025	128	8	4	4	0.1	0.01
Space Invaders	10M	256	0.00025	128	8	4	4	0.1	0.01
Seaquest	10M	256	0.00025	128	8	4	4	0.1	0.01
Beam Rider	10M	256	0.00025	128	8	4	4	0.1	0.01

Table 5: Uniform trained PPO hyperparameter configurations across all Atari environments. Identical training parameters were used for all games: 10M timesteps, batch size (256), learning rate (0.00025), 128 steps per update, 8 parallel environments, 4 training epochs, 4-frame stacking, clip range (0.1), and entropy coefficient (0.01).

Table 5 presents the consistent PPO <sup>6</sup> hyperparameters used for InfrectroRL attack across six Atari environments (Pong, Breakout, Qbert, Space Invaders, Seaquest, and Beam Rider). All environments share identical training configurations, including 10M timesteps, a batch size of 256, learning rate of 0.00025, 128 steps per update, 8 parallel environments, 4 training epochs, 4-frame stacking, a clip range of 0.1, and an entropy coefficient of 0.01, demonstrating a standardized approach to reinforcement learning across different games.

## Proof of Detectability Guarantees

In this appendix section we provide the full proof and a brief discussion of the main assumptions behind Theorem 2. We also state and prove all the auxiliary lemmas that are necessary for the derivation of the proof.

We start by arguing that the set of assumptions made for Theorem 2 is not excessively restrictive. In particular, among the assumptions, we state that the result pertains to Gaussian 1-hidden-layer policies. However, this assumption can be relaxed in two ways. First of all, the result is easily extendable to  $L$ -layers MLPs quite straightforwardly (see Lemma 5). This will impact coefficient  $B_j$  uniquely. Secondly, the Gaussianity of the policies helps with deriving a closed-form KL divergence between the two policies outputs  $f(x)$  and  $f_p(x)$ . This can be relaxed in favor of a general non-parametric assumption on the distribution form, but at the cost of losing a direct dependency between the difference in performance and the coefficient  $B_j$  (although one can derive closed-form KL divergence also for other classes of distributions, such as Binomial for instance). Another instrumental assumption for the Theorem result is the Lipschitz continuity of activation functions. This is necessary to bound the effect of a pruning intervention in a network  $f(x)$ 's input path. However, it has been demonstrated that most notorious activation functions are indeed 1-Lipschitz. This include ReLU, TanH and Softmax. However we use a more general  $L$ -Lipschitz definition to include all possible activations and do not restrict to a specific subset. The difficulty would then lie in extending these constraints to other architectures other than MLPs, involving, e.g., convolutions. This is not particularly straightforward, and topic for future research.

*Proof of Theorem 2.* Suppose we have two Gaussian, 1-hidden layer, policy networks  $\pi_1$  and  $\pi_2$ , and two time-dependent state-action transition densities  $p_1^t(s, a)$  and  $p_2^t(s, a)$ . Then we have:

$$\begin{aligned}
|J(\pi_1) - J(\pi_2)| &= \left| \sum_{t=0}^T \gamma^t \left[ \mathbb{E}_{a \sim \pi_1, s \sim p_1} [r(s_t, a_t)] - \mathbb{E}_{a \sim \pi_2, s \sim p_2} [r(s_t, a_t)] \right] \right| = \\
&= \left| \sum_{t=0}^T \gamma^t \left[ \int_a \int_s [p_1^t(s, a) - p_2^t(s, a)] r(s, a) ds da \right] \right| \leq \\
&\leq \sum_{t=0}^T \gamma^t \left[ \int_a \int_s |p_1^t(s, a) - p_2^t(s, a)| r(s, a) ds da \right]
\end{aligned}$$

<sup>6</sup><https://github.com/DLR-RM/rl-trained-agents/tree/cd35bde610f4045bf2e0731c8f4c88d22df8fc85>

Using the fact that rewards are bounded  $r : \mathcal{S} \times \mathcal{A} \rightarrow [r_{\min}, r_{\max}] \subset \mathbb{R}$ , then:

$$\begin{aligned}
|J(\pi_1) - J(\pi_2)| &\leq \sum_{t=0}^T \gamma^t \left[ \int_a \int_s |p_1^t(s, a) - p_2^t(s, a)| r(s, a) ds da \right] \leq \\
&\leq r_{\max} \sum_{t=0}^T \gamma^t \left[ \int_a \int_s |p_1^t(s, a) - p_2^t(s, a)| ds da \right] = \\
&= 2r_{\max} \sum_{t=0}^T \gamma^t \left[ D_{TV}(p_1^t(s, a) \| p_2^t(s, a)) \right]
\end{aligned}$$

Then, by Lemma 3 (see below) we have:

$$\begin{aligned}
|J(\pi_1) - J(\pi_2)| &\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ D_{TV}(p_1^t(s, a) \| p_2^t(s, a)) \right] \leq \\
&\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ D_{TV}(p_1^t(s) \| p_2^t(s)) + \mathbb{E}_{s \sim p} [D_{TV}(\pi_1(a|s) \| \pi_2(s, a))] \right]
\end{aligned}$$

Additionally, via Lemma 4 (see below) we obtain:

$$\begin{aligned}
|J(\pi_1) - J(\pi_2)| &\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ D_{TV}(p_1^t(s) \| p_2^t(s)) + \mathbb{E}_{s \sim p} [D_{TV}(\pi_1(a|s) \| \pi_2(s, a))] \right] \leq \\
&\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \mathbb{E}_{s \sim p} [D_{TV}(\pi_1(a|s) \| \pi_2(a|s))] \right]
\end{aligned}$$

From here, using Pinsker's inequality we derive:

$$\begin{aligned}
|J(\pi_1) - J(\pi_2)| &\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \mathbb{E}_{s \sim p} [D_{TV}(\pi_1(a|s) \| \pi_2(a|s))] \right] \leq \\
&\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \mathbb{E}_{s \sim p} [\sqrt{2D_{KL}(\pi_1(a|s) \| \pi_2(a|s))}] \right]
\end{aligned}$$

Considering that we have two policies  $\pi_1$  and  $\pi_2$  that are 1-hidden layer MLPs, outputting a Gaussian distribution, by Lemma 5 below we obtain:

$$\begin{aligned}
|J(\pi_1) - J(\pi_2)| &\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \mathbb{E}_{s \sim p} [\sqrt{2D_{KL}(\pi_1(a|s) \| \pi_2(a|s))}] \right] \leq \\
&\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \mathbb{E}_{s \sim p} \left[ \sqrt{2 \frac{(B_j s_j)^2}{2\sigma_f^2}} \right] \right] = \\
&= 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \frac{B_j \mathbb{E}_{s \sim p} [|s_j|]}{\sigma_f} \right]
\end{aligned}$$

Finally, using the normalized inputs assumption by which  $s_j \in [0, 1]$  then we have:

$$\begin{aligned}
|J(\pi_1) - J(\pi_2)| &\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \frac{B_j \mathbb{E}_{s \sim p} [|s_j|]}{\sigma_f} \right] \leq \\
&\leq 2r_{\max} \sum_{t=0}^T \gamma^t \left[ t\delta + \frac{B_j}{\sigma_f} \right] \leq \\
&\leq 2r_{\max} \sum_{t=0}^{\infty} \gamma^t \left[ t\delta + \frac{B_j}{\sigma_f} \right] \leq \\
&\leq 2r_{\max} \left[ \frac{\gamma\delta}{(1-\gamma)^2} + \frac{B_j}{\sigma_f(1-\gamma)} \right],
\end{aligned}$$

where  $\frac{1}{2} \mathbb{E}_{s' \sim p} [D_{TV}(p_1(s|s') \| p_2(s|s'))] \leq \delta$  according to assumptions.  $\square$

After stating the full proof of Theorem 2, we proceed here below by providing statements of all the auxiliary lemmas instrumental to the theorem's proof, together with their own standalone proofs. For almost all the lemmas and proofs, we will be using the Total-Variation Distance on continuous probability spaces, which is defined for a continuous random variable  $s \in \mathcal{S}$  and two density functions  $p, q$  defined on the same probability space as:

$$D_{TV}(p(s)||q(s)) = \frac{1}{2} \int_{s \in \mathcal{S}} |p(s) - q(s)| ds .$$

**Lemma 3** (Joint Probability  $D_{TV}$  Decomposition). *Consider two, time-dependent, joint state-action visitation probabilities  $p_1^t(s, a)$  and  $p_2^t(s, a)$ , and their Total Variation Distance:*

$$D_{TV}(p_1^t(s, a)||p_2^t(s, a)) = \frac{1}{2} \int_s \int_a |p_1^t(s, a) - p_2^t(s, a)| ds da .$$

We can decompose this quantity into:

$$D_{TV}(p_1^t(s, a)||p_2^t(s, a)) \leq D_{TV}(p_1^t(s)||p_2^t(s)) + \mathbb{E}_{s \sim p}[D_{TV}(\pi_1(a|s)||\pi_2(a|s))] .$$

*Proof of Lemma 3.* Let us first define the recursive one-step decomposition of the joint probability  $p_i^t(s, a)$ :

$$p_i^t(s, a) = \int_s p_i(s'|s, a) \pi_i(a|s) p_i^{t-1}(s) ds .$$

Then, we have that the difference:

$$\begin{aligned} p_1^{t+1}(s, a) - p_2^{t+1}(s, a) &= \int_s [p_1(s'|s, a) \pi_1(a|s) p_1^t(s) - p_2(s'|s, a) \pi_2(a|s) p_2^t(s)] ds = \\ &= \int_s \left\{ p_1(s'|s, a) [\pi_1 - \pi_2] p_1^t(s) + p_2(s'|s, a) \pi_2 [p_1^t(s) - p_2^t(s)] \right\} ds \end{aligned}$$

The Total Variation Distance between the two densities is then defined as:

$$\begin{aligned} D_{TV}(p_1^{t+1}(s, a)||p_2^{t+1}(s, a)) &= \frac{1}{2} \int_s \int_a |p_1^{t+1}(s, a) - p_2^{t+1}(s, a)| da ds \leq \\ &\leq \frac{1}{2} \int_s \int_a |p_1(s'|s, a) [\pi_1 - \pi_2] p_1^t(s)| ds da + \frac{1}{2} \int_s \int_a |p_2(s'|s, a) \pi_2 [p_1^t(s) - p_2^t(s)]| ds da \end{aligned}$$

For the first term of the sum above, we can derive:

$$\begin{aligned} &\frac{1}{2} \int_s \int_a |p_1(s'|s, a) [\pi_1 - \pi_2] p_1^t(s)| ds da \leq \\ &\leq \frac{1}{2} \int_s \int_a |\pi_1 - \pi_2| p_1^t(s) ds da \leq \frac{1}{2} \int_s \int_a |\pi_1 - \pi_2| ds da \leq \\ &\int_s D_{TV}(\pi_1||\pi_2) ds = \mathbb{E}_{s \sim p}[D_{TV}(\pi_1||\pi_2)] \end{aligned}$$

For the second term instead we have:

$$\begin{aligned} &\frac{1}{2} \int_s \int_a |p_2(s'|s, a) \pi_2 [p_1^t(s) - p_2^t(s)]| ds da \leq \frac{1}{2} \int_s \int_a |p_1^t(s) - p_2^t(s)| ds da \leq \\ &\frac{1}{2} \int_s |p_1^t(s) - p_2^t(s)| ds = D_{TV}(p_1^t(s)||p_2^t(s)) \end{aligned}$$

Putting the two together, we eventually obtain:

$$D_{TV}(p_1^{t+1}(s, a)||p_2^{t+1}(s, a)) \leq D_{TV}(p_1^t(s)||p_2^t(s)) + \mathbb{E}_{s \sim p}[D_{TV}(\pi_1(a|s)||\pi_2(a|s))] .$$

□

**Lemma 4** (Time-Dependent  $D_{TV}$  Bound). *Suppose we assume same initial state distributions  $p_1^0(s) = p_2^0(s)$  and we define the following supremum quantity  $\sup_t \mathbb{E}_{\tilde{s} \sim p}[D_{TV}(p_1(s|\tilde{s})||p_2(s|\tilde{s}))] \leq \delta$ . Then we can obtain the following bound:*

$$D_{TV}(p_1^t(s)||p_2^t(s)) \leq t\delta .$$

*Proof of Lemma 4.* Let us start by defining the one-step recursive decomposition:

$$p_1^t(s) = \int_{\tilde{s}} p_1(s|\tilde{s}) p_1^{t-1}(\tilde{s}) d\tilde{s} .$$

Then we have:

$$\begin{aligned} D_{TV}(p_1^t(s) \| p_2^t(s)) &= \frac{1}{2} \int_s |p_1^t(s) - p_2^t(s)| ds = \\ &= \frac{1}{2} \int_{\tilde{s}} \int_{\tilde{s}} |p_1(s|\tilde{s}) p_1^{t-1}(\tilde{s}) - p_2(s|\tilde{s}) p_2^{t-1}(\tilde{s})| ds d\tilde{s} \leq \\ &\leq \frac{1}{2} \int_s \left[ \int_{\tilde{s}} p_1^{t-1}(\tilde{s}) |p_1(s|\tilde{s}) - p_2(s|\tilde{s})| d\tilde{s} + \int_{\tilde{s}} p_2^{t-1}(\tilde{s}) |p_1^{t-1}(s) - p_2^{t-1}(s)| d\tilde{s} \right] ds = \\ &= \frac{1}{2} \int_{\tilde{s}} p_1^{t-1}(\tilde{s}) \left[ \int_s |p_1(s|\tilde{s}) - p_2(s|\tilde{s})| ds \right] d\tilde{s} + \frac{1}{2} \int_{\tilde{s}} |p_1^{t-1}(\tilde{s}) - p_2^{t-1}(\tilde{s})| d\tilde{s} = \\ &= \mathbb{E}_{\tilde{s} \sim p^{t-1}} \left[ \frac{1}{2} \int_s |p_1(s|\tilde{s}) - p_2(s|\tilde{s})| ds \right] + D_{TV}(p_1^{t-1}(s) \| p_2^{t-1}(s)) \leq \\ &\leq \delta + D_{TV}(p_1^{t-1}(s) \| p_2^{t-1}(s)) , \end{aligned}$$

By recursion over time  $t$ , we can compactly rewrite:

$$D_{TV}(p_1^t(s) \| p_2^t(s)) \leq t\delta + D_{TV}(p_1^0(s) \| p_2^0(s)) ,$$

where by assumption  $D_{TV}(p_1^0(s) \| p_2^0(s)) = 0$ , leaving:

$$D_{TV}(p_1^t(s) \| p_2^t(s)) \leq t\delta .$$

□

Finally here below we prove the lemma regarding the masking of Multi-Layer Perceptron networks.

**Lemma 5** (Path Masking Intervention on MLPs). *Suppose we have a 1-hidden layer MLP that outputs  $y \sim \mathcal{N}(f(x), \sigma_f^2)$  from inputs  $(x_1, \dots, x_d)$ , where  $\sigma_f^2 \in \mathbb{R}^+$  and:*

$$f(x) = \sum_{i=1}^H W_{2,i} \phi(z_i) + b_2 , \quad \text{where } z_i = \sum_{k=1}^d W_{1,ik} x_k + b_{1,i} ,$$

and  $\phi : \mathbb{R} \rightarrow (a, b) \subseteq \mathbb{R}$  is a  $L_\phi$ -Lipschitz continuous activation function. Then defining as  $f_p(x)$  the “pruned” version of  $f(x)$ , where input  $x_j$  is masked out via the intervention  $W_{1,j} = 0$ , we have that:

$$D_{KL}(\mathcal{N}(f(x), \sigma_f^2) \| \mathcal{N}(f_p(x), \sigma_f^2)) \leq \frac{(B_j |x_j|)^2}{2\sigma_f^2} ,$$

where  $B_j = L_\phi \sum_{i=1}^H |W_{1,ij} W_{2,i}|$ .

*Proof of Proposition 5.* The “pruned” version  $f_p(x)$  of  $f(x)$  can be written in the following form:

$$f_p(x) = \sum_{i=1}^H W_{2,i} \phi(z'_i) + b_2 , \quad \text{where } z'_i = \sum_{k=1, k \neq j}^d W_{1,ik} x_k + b_{1,i} .$$

Then, it is easy to obtain that  $z'_i = z_i - W_{1,ij} x_j$ , so that the difference becomes:

$$f(x) - f_p(x) = \sum_{i=1}^H W_{2,i} [\phi(z_i) - \phi(z_i - W_{1,ij} x_j)] .$$

Using the  $L_\phi$ -Lipschitz property of  $\phi(\cdot)$  we obtain:

$$\begin{aligned} |f(x) - f_p(x)| &= \left| \sum_{i=1}^H W_{2,i} [\phi(z_i) - \phi(z_i - W_{1,ij} x_j)] \right| \leq \sum_{i=1}^H |W_{2,i}| |\phi(z_i) - \phi(z_i - W_{1,ij} x_j)| \leq \\ &\leq \sum_{i=1}^H |W_{2,i}| L_\phi |W_{1,ij} x_j| = |x_j| L_\phi \sum_{i=1}^H |W_{2,i} W_{1,ij}| = B_j |x_j| . \end{aligned}$$

Finally this implies that:

$$D_{KL}(\mathcal{N}(f(x), \sigma_f^2) \parallel \mathcal{N}(f_p(x), \sigma_f^2)) = \frac{(f(x) - f_p(x))^2}{2\sigma_f^2} \leq \frac{(B_j|x_j|)^2}{2\sigma_f^2},$$

where  $D_{KL}$  is the KL divergence, defined as  $D_{KL}(p \parallel q) = \int_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx$ .

□

## Ablation Study

We conduct ablation studies on InfrectroRL’s hyperparameters: threshold ( $\lambda$ ), amplification factor ( $\gamma$ ), trigger size, and target label. For control, each hyperparameter is set to suboptimal values to observe episodic return deviations.

**Impact of  $\gamma$**  We observe a consistent negative correlation between  $\gamma$  and average reward (with  $\lambda$ , trigger size, and target label fixed). This trend is expected: higher  $\gamma$  amplifies neurons, increasing the target action probability and consequently lowering episodic reward.

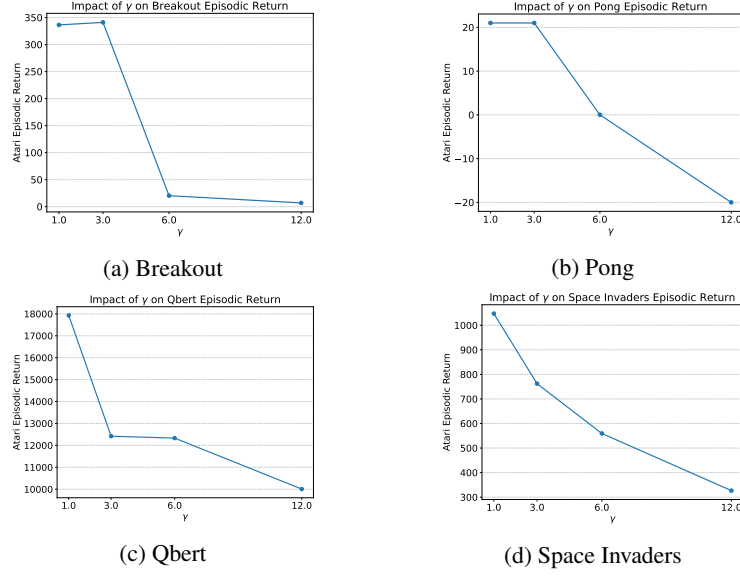


Figure 4: Ablation study in InfrectroRL for varying  $\gamma$  in all 4 games. We notice that varying  $\gamma$  severely reduces the episodic return, as expected from our experimentations.  $\gamma$  is known as the amplification factor that traverses through the policy network to induce a malicious adversary-chosen action.

**Impact of  $\lambda$**  We notice that our attack maintains effective levels of episodic return regardless of the variation in  $\lambda$ . The reason is because the backdoor path is always activated for backdoored inputs where  $\lambda > 0$ .

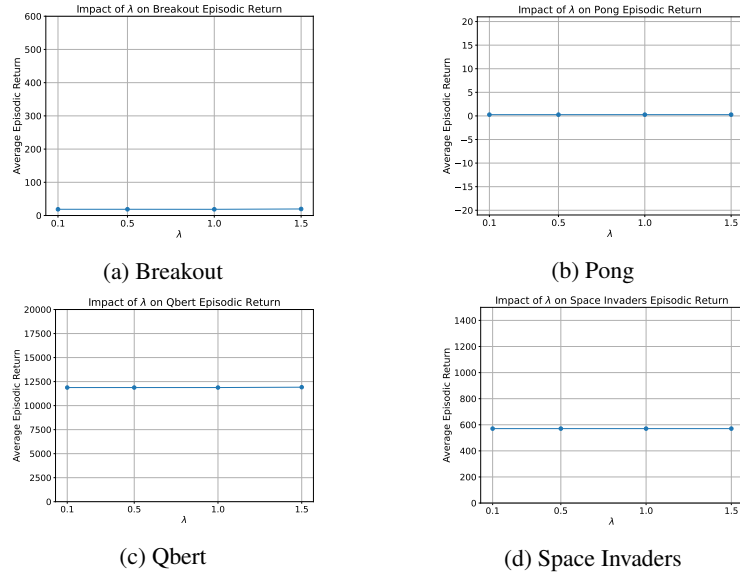


Figure 5: Ablation study in InfrectroRL for varying  $\lambda$  in all 4 games. We notice that varying  $\lambda$  has no effect on the episodic return. The reason is that the backdoor path crafted by InfrectroRL is always activated for backdoored inputs when  $\lambda > 0$ .



**Impact of Trigger Size** We notice that our attack consistently reduces the episodic reward regardless of the trigger size, making our attack resilient to change in appropriate trigger size.

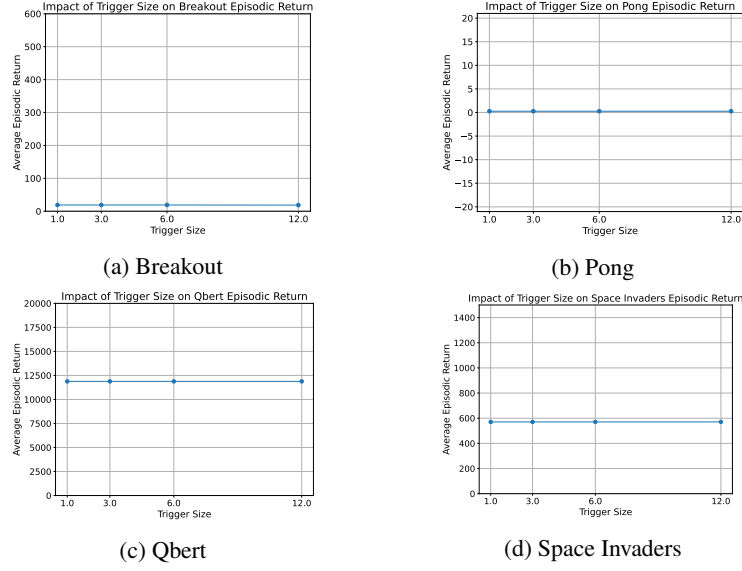


Figure 6: Ablation study in InfrectroRL for varying the trigger size from 1 to 12 in all 4 games. We notice that InfrectroRL is highly effective regardless of the trigger size.

**Impact of Target Label** Given that our agent is a decision-making algorithm, we assume that the target actions can also affect the overall episodic return. As we set the hyperparameters, we set the remaining hyperparameters to a static values. We notice that all target label actions amount to similar levels of performance in the average episodic return. We assume this is primarily because of the repetitive actions made by the agent that cause it to reach a corner (or stay in the exact same position) in the environment.

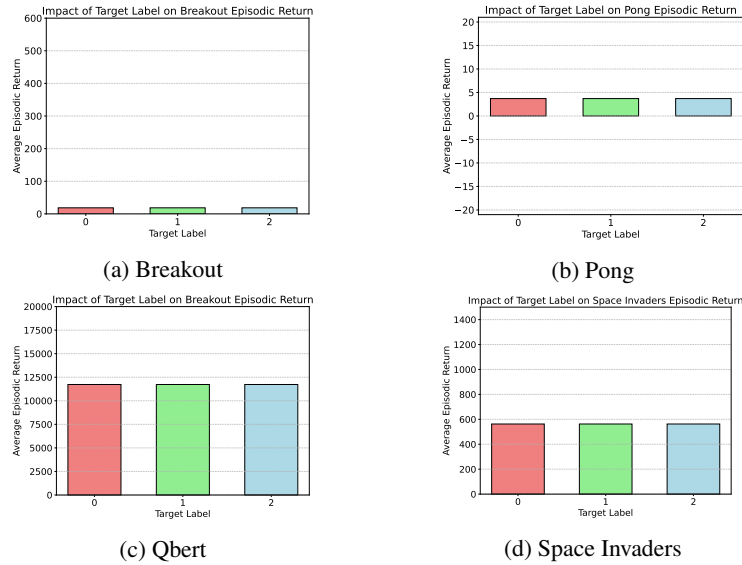


Figure 7: Ablation study in InfrectroRL for varying the target label for all 4 games. We notice that the backdoor effects of InfrectroRL on the episodic return is resilient regardless of the target label.