# Adversarial Attacks through Value-Guided Transition Modeling in Deep Reinforcement Learning

Thomas O'Cuilleanain<sup>1</sup>, Juan Cardenas-Cartagena<sup>1</sup>, and Matthia Sabatelli\*<sup>1</sup>

<sup>1</sup>University of Groningen

1m.sabatelli@rug.nl

### Abstract

001

003

004

005

006

007

008

010

011

012

014

017

018

019

020

021

022

023

024

025

026

028

029

030

031

032

033

034

035

Efficient adversarial attacks on deep reinforcement learning agents rely on identifying critical states. Prior work uses learned transition models with environment-specific metrics to predict and lure the victim agent to such states. We propose a valueguided attack that integrates the victim policy's value function as an environment-agnostic metric into both transition model training and state evaluation. From our preliminary results in the Pong environment from the Arcade Learning Environment, our method achieves comparable performance degradation to prior work while requiring roughly half as many attacks.

#### Background 1 015

Reinforcement Learning: Reinforcement Learning (RL) environments are modeled as Markov Decision Processes (MDP) [1], defined by

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, \mathcal{R}, \gamma),$$

where S is the state space, A the action space,  $p(s_{t+1} \mid s_t, a_t)$  the transition function,  $\mathcal{R}(s_t, a_t, s_{t+1})$ the reward function, and  $\gamma \in [0,1]$  the discount factor. At each time-step t, the agent observes the current state  $s_t \in \mathcal{S}$ , selects an action  $a_t \in \mathcal{A}$  according to its policy  $\pi$ , transitions to a next state  $s_{t+1} \sim p(\cdot \mid s_t, a_t)$  and receives a scalar reward  $r_t = \mathcal{R}(s_t, a_t, s_{t+1})$ . The agent's behavior is governed by its policy  $\pi(a_t \mid s_t)$ , which defines a probability distribution over actions given the current state. The state-value function represents the expected cumulative discounted return when starting from  $s_t$  and following  $\pi$  thereafter:

$$V^{\pi}(s) = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t = s \right]. \tag{1}$$

In this work, we leverage the value function,  $V^{\pi}(s)$ , as an environment-agnostic metric for identifying critical states to which we lure the victim policy.

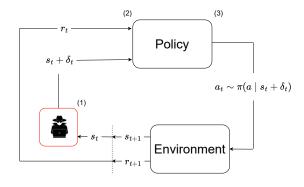


Figure 1. General schema of test-time, state-based adversarial attacks in DRL: (1) the adversary intercepts the true observation and injects a perturbation  $\delta_t$  into the state  $s_t$  (2) which is passed to the agent. (3) The agent then samples an action  $a_t \sim \pi(a \mid s_t + \delta_t)$  which is executed in the environment.

Adversarial Attacks: Test-time adversarial at- 037 tacks on Deep Reinforcement Learning (DRL) agents aim to manipulate the behavior of a trained victim policy to achieve an ulterior objective (e.g. performance degradation). These attacks exploit the sensitivity of trained policies by perturbing their input in order to induce the policy into taking (targeted) suboptimal actions during evaluation. We denote the perturbation added to  $s_t$  as  $\delta_t$ . Figure 1 illustrates this framework.

Sun et al. [2] propose Critical Point Attack (CPA) which identifies critical states to make adversarial attacks more efficient by learning a parametrized transition model  $f_{\theta}: \mathcal{S} \times \mathcal{A} \to S$  of the environment, where  $\theta$  are the parameters. This model follows the same architecture proposed by Oh et al. [3]. Given a dataset of N collected trajectories from the victim agent  $D = \left\{ \left( (s_0^{(i)}, a_0^{(i)}), \ldots, (s_{T_i}^{(i)}, a_{T_i}^{(i)}) \right) \right\}_{i=0}^N$ , and a prediction horizon with length K,  $\theta$  is optimised using the loss,

$$SE^{(i)}(t;\theta) = \|\hat{s}_t^{(i)}(\theta) - s_t^{(i)}\|_2^2,$$
 (2) 057

046

047

049

053

054

055

056

$$L_{\text{CPA}}(\theta) = \frac{1}{2K} \sum_{i,t} \sum_{k=1}^{K} \text{SE}^{(i)}(t+k;\theta),$$
 (3) of

<sup>\*</sup>Corresponding Author.

<sup>&</sup>lt;sup>1</sup>Here "transition model" denotes a learned environment dynamics predictor, not the MDP transition function.

059

060

061

062

063

064

065

067

068

069

070

071

072

075

076

078

079

080

081

083

084

085

086

087

088

090

091

092

093

094

095

098

099

100

103

104

105

106

107

108

113

114

117

118

119

126

133

134

135

137

141

143

144

where  $\hat{s}_t^{(i)}(\theta) := f_{\theta}(\hat{s}_{t-1,\theta}^{(i)}, a_{t-1}^{(i)}), t > 0$ . For simplicity, we omit  $\theta$  in the  $\hat{s}$  notation unless required explicitly. Using  $f_{\theta}$ , given a rollout horizon length M, and the victim policy  $\pi$ , the adversary predicts the baseline state  $\hat{s}_{t+M}^{\pi}$  starting from  $s_t$  by following the victim policy on predicted next states recursively. Next, given an attack horizon of length  $K \leq M$ , the adversary predicts all possible subsequent states by enumerating all possible action sequences of length K. We denote a specific sequence as  $\mathbf{a}_{t:t+K} \triangleq (a_i)_{i=t}^{t+K} \in \mathcal{A}^K$ , where  $\mathcal{A}^K$  is the Cartesian product of A with itself K times. If M > K, then for each such predicted state the adversary continues the rollout using the victim policy for the remaining M-K steps. This results in  $|\mathcal{A}^K|$ predicted final states. Using an environment specific divergence function  $T: \mathcal{S} \to \mathbb{R}$ , the adversary finds the final action sequence from the starting state which maximises the Danger Awareness Metric,

$$DAM_T(\mathbf{a}_{t:t+K}) = \left| T(\hat{s}_{t+M}^{\mathbf{a}_{t:t+K}}) - T(s_{t+M}^{\pi}) \right|, \quad (4)$$

$$\mathbf{a}_{t:t+K}^* = \underset{\mathbf{a}_{t:t+K} \in \mathcal{A}^K}{\operatorname{arg\,max}} \operatorname{DAM}_T(\mathbf{a}_{t:t+K}), \quad (5)$$

where  $\hat{s}_{t+M}^{\mathbf{a}_{t:t+K}}$  denotes the predicted state  $\hat{s}$  at time step t+M following action sequence  $\mathbf{a}_{t:t+K}$  for K-steps and the victim policy  $\pi$  for M-K-steps thereafter. If, for any given final state, this metric surpasses a threshold  $\Delta > 0$ , the victim policy is fooled into following the associated action sequence by adding carefully crafted perturbations. These perturbations are computed using the Carlini & Wagner (C&W) attack [4]. The authors state that it is necessary for T to have an environment-specific definition to accurately reflect the potential danger associated with a predicted state. In the Pong and Breakout environments from the Arcade Learning Environment [5], the authors turn to predicting the RAM state representation of subsequent states. Given the RAM state s, the authors define  $T(s) = d(s) \cdot p(s)$ , where  $d: \mathcal{S} \to \mathbb{R}$  is the Euclidean distance between the ball and the paddle and  $p: \mathcal{S} \to \{0,1\}$  is equal to 1 if the ball has been dropped and 0 otherwise.

### Our Method $\mathbf{2}$

Rationale: The proposed attack employs an environment-agnostic T function based on the statevalue function  $V^{\pi}$ . In this context, the adversary finds the action sequence  $\mathbf{a}_{t:t+K}^*$  as follows,

$$\mathbf{a}_{t:t+K}^* = \underset{\mathbf{a}_{t:t+K} \in \mathcal{A}^K}{\operatorname{arg \, max}} \operatorname{DAM}_{V^{\pi}}(\mathbf{a}_{t:t+K}). \tag{6}$$

From our early experiments using  $L_{CPA}$  to optimize  $\theta$ , we observed that in the Pong environment, the final reconstruction loss was minimal, yet the difference in state-value estimates between genuine and predicted states remained large, ultimately hindering the effectiveness of the attack. We hypothesize

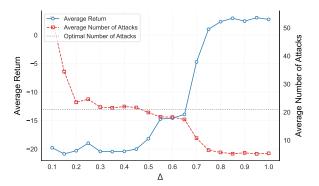


Figure 2. Average return and number of attacks vs  $\Delta$ against trained A2C victim policy using value-guided attack with K = M = 1 in Pong averaged over 100 episodes.

that this occurs because, although the reconstruction error is low, the ball in these environments occupies only a few pixels in the observation space. Consequently, the reconstruction loss provides a very weak learning signal for its precise position. These inaccuracies in the ball's position lead to substantial discrepancies in the policy's value estimates, which are highly sensitive to ball position. Therefore, we propose a new loss to optimize  $\theta$ ,

$$\mathrm{SE}_{V^{\pi}}^{(i)}(t;\theta) = \left\|V^{\pi}\big(\hat{s}_t^{(i)}(\theta)\big) - V^{\pi}(s_t^{(i)})\right\|_2^2, \quad (7) \quad \ \ \, \text{120}$$

$$L_V(\theta) = \frac{1}{2K} \sum_{i,t} \sum_{k=1}^K SE_{V^{\pi}}^{(i)}(t+k;\theta),$$
 (8) 121

$$L(\theta) = \alpha L_{\text{CPA}}(\theta) + (1 - \alpha)L_V(\theta), \qquad (9) \quad {}_{122}$$

where  $\alpha \in [0,1]$  allows tuning for different tasks. 123 This guides the reconstruction to produce subsequent states that are not only visually accurate but aligned with  $V^{\pi}$ , thereby emphasizing task-relevant features, such as the position of the ball.

Preliminary Findings: In our experiments, we optimize  $\theta$  using L with  $\alpha = 0.1$ . We evaluate this attack against a trained victim A2C agent in the Pong environment. The lowest return of -21 is achieved after a minimum of 21 attacks, as the agent must drop the ball 21 times. Averaging over 100 episodes, with K = M = 1, and threshold  $\Delta = 0.35$ the attack achieves an average return of -20.45 in 21.55 attacks (see Figure 2). Compared to the results in Sun et al. [2], the attack achieves comparable performance degradation in around half the number of attacks (K = M = 2). We plan to test our environment-agnostic adversarial attack in environments in which constructing an environmentspecific T function is infeasible; yet estimating the victim policy's value function remains tractable for discrete action spaces.

## References

- 146 [1] M. L. Puterman. "Markov decision processes". 147 In: *Handbooks in operations research and man-*148 agement science 2 (1990), pp. 331–434.
- 149 [2] J. Sun, T. Zhang, X. Xie, L. Ma, Y. Zheng,
  150 K. Chen, and Y. Liu. "Stealthy and efficient
  151 adversarial attacks against deep reinforcement
  152 learning". In: Proceedings of the AAAI confer153 ence on artificial intelligence. Vol. 34. 04. 2020,
  154 pp. 5883–5891.
- [3] J. Oh, X. Guo, H. Lee, R. L. Lewis, and S. Singh.
   "Action-conditional video prediction using deep
   networks in atari games". In: Advances in neural information processing systems 28 (2015).
- 159 [4] N. Carlini and D. Wagner. "Towards evaluating 160 the robustness of neural networks". In: 2017 161 ieee symposium on security and privacy (sp). 162 Ieee. 2017, pp. 39–57.
- [5] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. "The Arcade Learning Environment:
   An Evaluation Platform for General Agents".

   In: Journal of Artificial Intelligence Research
   47 (June 2013), pp. 253–279.