

Latent Action Discretization in DiLA for Robust Multimodal Language-Conditioned Policies Under Visual Perception Noise on CALVIN

Assignee Research

June 12, 2026

Abstract

This paper investigates the resilience and robustness of Deep Reinforcement Learning (DRL) policies to adversarial perturbations in the state space. We first present an approach for the disentanglement of vulnerabilities caused by representation learning of DRL agents from those that stem from the sensitivity of the DRL policies to distributional shifts in state transitions. Building on this approach, we propose two RL-based techniques for quantitative benchmarking of adversarial resilience and robustness in DRL policies against perturbations of state transitions. We demonstrate the feasibility

1 Introduction

This paper examines: RL-Based Method for Benchmarking the Adversarial Resilience and Robustness of Deep Reinforcement Learning Policies. Research question: To what extent does the latent action discretization strategy in DiLA improve robustness against visual perception noise in multimodal language-conditioned policies evaluated on CALVIN?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.4/10.

3 Results

15 papers retrieved. 12 claims extracted; 9 independently verified. Quality review score: 7.4/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
For A2C and PPO2 targets, the state-action value function $Q^*(st, a)$ is calculated using the formula $Q^*(st, a) = r(st, a)$	✓	0.23
All three adversarial policies (targeting DQN, A2C, and PPO2) converge to the same optima during training.	×	0.14
The adversary targeting the DQN policy requires a higher number of training steps to achieve convergence compared to adv	✓	0.16
At convergence, the mean-per-100 episodes of the minimum number of perturbations for the A2C target is 7.69.	✓	0.18
At convergence, the mean-per-100 episodes of the minimum number of perturbations for the PPO2 target is 7.49.	✓	0.20
At convergence, the mean-per-100 episodes of the minimum number of perturbations for the DQN target is 7.13.	✓	0.20
During test-time, the DQN policy requires 6.95 perturbations to incur an adversarial regret of 491.15.	✓	0.19
During test-time, the PPO2 policy requires 7.72 perturbations to incur an adversarial regret of 490.47.	✓	0.18
During test-time, the A2C policy requires 8.71 perturbations to incur an adversarial regret of 488.16.	✓	0.17
The DQN policy has the lowest adversarial resilience among the three tested policies (DQN, PPO2, A2C).	×	0.15
The A2C policy is the most resilient to state-space perturbation attacks among the three tested policies.	✓	0.16
For all three adversarial policies, the initial timesteps of an episode are the most frequently perturbed.	×	0.10

References

- <http://arxiv.org/abs/2509.19212v1>

- <http://arxiv.org/abs/1906.01110v1>
- <http://arxiv.org/abs/2605.15725v1>