

Robust Deep Reinforcement Learning against Adversarial Behavior Manipulation

ICLR2026

Shojiro Yamabe¹, Kazuto Fukuchi^{2,3}, Jun Sakuma^{1,3}

¹ Institute of Science Tokyo, ² University of Tsukuba, ³ RIKEN AIP

- Deep Reinforcement Learning (DRL) agents are vulnerable to adversarial attack, like classification model
- In recent years, applications have expanded to mission-critical tasks
 - Alignment of LLM [1]: Safety alignment to avoid harmful outputs
 - Autonomous driving [2]
 - Factory automation [3]
- As these developments progress, the importance of reinforcement learning security has increased

[1] B Ravi Kiran, et al. Deep reinforcement learning for autonomous driving: A survey. IEEE Transactions on Intelligent Transportation Systems, 2021.

[2] Long Ouyang, et al. Training language models to follow instructions with human feedback. NeurIPS, Vol. 35, pp. 27730–27744, 2022.

[3] Jonas Degraeve, et al. Magnetic control of tokamak plasmas through deep reinforcement learning. Nature, 2022.

- Many studies have explored adversarial attacks and defenses in reinforcement learning
- However, most prior work focuses on reward-minimization attacks [4,5,6,7]
 - which aim to degrade the victim's performance by reducing its reward.
- But attackers may have other goals.
 - E.g., they may force a self-driving car to detour through a specific store
 - E.g., they may manipulate a recommender system to promote beneficial products.
- These goals are defined **independently** of the victim's reward.

 We therefore focus on the **victim's behavior** itself, rather than reward

[4] Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. *ICLR*, 2017.

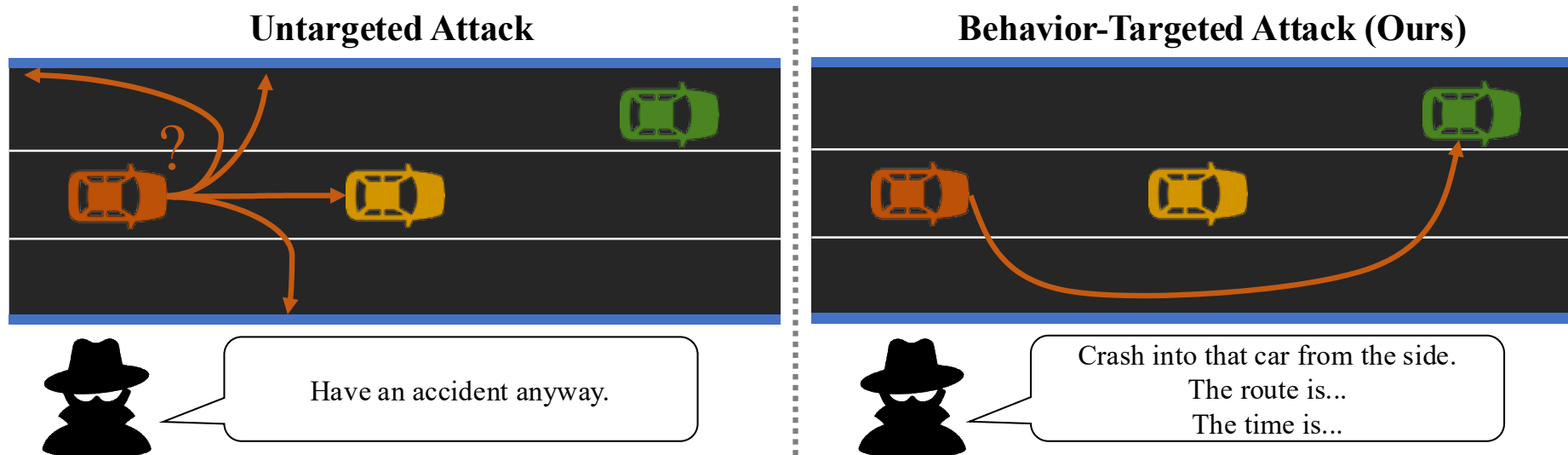
[5] Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommanan, and Girish Chowdhary. Robust deep reinforcement learning with adversarial attacks. IFAAMAS, 2018.

[6] Zhang, Huan, et al. "ROBUST REINFORCEMENT LEARNING ON STATE OBSERVATIONS WITH LEARNED OPTIMAL ADVERSARY.?", *ICLR 2021*.

[7] Sun, Yanchao, et al. "Who Is the Strongest Enemy? Towards Optimal and Efficient Evasion Attacks in Deep RL." , *ICLR 2022*

Behavior-targeted attack

- In this work, we propose the **behavior-targeted attack** and its **countermeasure** by intervening the agent's observations
 - Intervening observation amounts to attacking an autonomous vehicle's LiDAR or cameras
 - Behavior-targeted attacks aim to force the victim to behave specified by the adversary
 - A planned attack in which the location and timing of the accident are specified
- This poses a threat of far more sophisticated and stealthy attacks than untargeted attack whose only goal is to only degrade performance



- We first integrate **imitation learning (IL)** to train adversarial policy
 - IL learns a policy by observing expert demonstrations and mimicking them
 - E.g., learning safe and efficient driving skills from data collected from expert drivers [8]
- Our approach overcomes the limitations of previous studies on behavior-targeted attack [9,10]
 - It relies on environmental heuristics (e.g., it can only be applied to autonomous vehicles).
 - It assume white-box access to the victim's policy

Contribution

Propose Behavior-Imitation Attack (BIA) which is applicable under **Not relies on environmental heuristics** and **Limited access to the victim**

[8] Hawke, Jeffrey, et al. "Urban driving with conditional imitation learning." *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020.

[9] Leonard Hussenot, et al. Copycat: Taking control of neural policies with constant attacks. In *AAMAS*, 2020.

[10] Adith Bloor, et al. Attacking vision-based perception in end-to-end autonomous driving models. *Journal of Systems Architecture*, Vol. 110, p. 101766, 2020.

- We formulate the defender's objective with attacker's gain R_{tgt}

$$\arg \min_{\pi} J_{\text{def}}(\pi) = \underbrace{-J_{\text{RL}}(\pi)}_{\substack{\text{Standard RL objective} \\ \text{(Reward Maximization)}}} + \lambda \underbrace{\left(\max_{\nu} \mathbb{E}_{\pi \circ \nu}^M [R_{\text{tgt}}(s, a)] - \mathbb{E}_{\pi}^M [R_{\text{tgt}}(s, a)] \right)}_{\text{Worst-case attacker's gain}},$$

- Our theoretical analysis reveals two important insights:
 - **suppressing the sensitivity of the policy's action outputs to state changes enhances robustness against attacks → Policy Smoothing is effective!**
 - achieving lower sensitivity during the early stages of a trajectory significantly improves overall robustness

Contribution

Propose the first robust training method against behavior-targeted attack, **Time-Discounted Robust Training (TDRT)**

- Environment: Meta-World [11]
- The further to the right, the stronger the assumptions imposed on the attacker
- Our method (BIA) achieves strong attack performance even under these stronger attacker assumptions

Adv Task	Target Reward	Attack Rewards (\uparrow)					Random (no-box, no knowledge)
		Adversary with full knowledge \leftarrow Targeted PGD (white-box, target policy)	Rew Max (PA-AD) (white-box, reward function)	Rew Max (SA-RL) (black-box, reward function)	BIA-ILfD (ours) (black-box, demonstrations)	BIA-ILfO (ours) (no-box, demonstrations)	
window-close	4543 \pm 39	1666 \pm 936	4255 \pm 300	4505 \pm 65	3962 \pm 666	4036 \pm 510	947 \pm 529
window-open	4508 \pm 121	515 \pm 651	493 \pm 562	506 \pm 444	566 \pm 523	557 \pm 679	322 \pm 261
drawer-close	4868 \pm 6	2891 \pm 150	3768 \pm 1733	4658 \pm 747	4760 \pm 640	4626 \pm 791	1069 \pm 1585
drawer-open	4713 \pm 16	953 \pm 450	1607 \pm 355	1499 \pm 536	1556 \pm 607	1445 \pm 610	841 \pm 357
faucet-close	4754 \pm 15	1092 \pm 192	1241 \pm 501	3409 \pm 652	3316 \pm 648	3041 \pm 502	897 \pm 171
faucet-open	4544 \pm 800	2541 \pm 86	1420 \pm 85	1448 \pm 64	3031 \pm 1493	2718 \pm 1293	1372 \pm 81
handle-press-side	4546 \pm 721	1994 \pm 1225	4726 \pm 175	4625 \pm 175	4631 \pm 408	4627 \pm 586	1865 \pm 1340
handle-pull-side	4442 \pm 732	2198 \pm 1524	2065 \pm 1501	3617 \pm 1363	4268 \pm 740	4193 \pm 517	1426 \pm 1617
door-lock	3845 \pm 79	640 \pm 664	763 \pm 768	1937 \pm 1186	2043 \pm 1229	1906 \pm 1045	589 \pm 494
door-unlock	4690 \pm 33	531 \pm 61	3295 \pm 1111	3421 \pm 974	3336 \pm 932	3123 \pm 1123	391 \pm 59

- Adversarial training, which is effective against reward-minimization attacks, is ineffective here.
- Defense methods that aim to smooth the policy achieve strong robustness.

Task	Adversarial Training			Policy Smoothing	
	PPO (No defense)	ATLA-PPO (AdvTraining)	PA-ATLA-PPO (AdvTraining)	SA-PPO Smoothing (w/o time-discounting)	TDRT-PPO (ours) Smoothing (w/ time-discounting)
window-close	4505 ± 65	4270 ± 188	4041 ± 96	485 ± 61	482 ± 3
window-open	566 ± 523	586 ± 649	671 ± 589	272 ± 37	254 ± 214
drawer-close	4760 ± 640	4858 ± 6	4868 ± 3	4 ± 2	4860 ± 4
drawer-open	1556 ± 607	1158 ± 1026	954 ± 219	403 ± 49	378 ± 10
faucet-close	3409 ± 652	4108 ± 790	4012 ± 123	1559 ± 406	1789 ± 610
faucet-open	3031 ± 1493	4383 ± 449	2358 ± 976	1763 ± 255	1942 ± 261
handle-press-side	4726 ± 175	4302 ± 799	3318 ± 1539	1888 ± 1169	1928 ± 736
handle-pull-side	4268 ± 740	532 ± 534	512 ± 982	10 ± 1	7 ± 1
door-lock	2043 ± 1229	1020 ± 805	992 ± 19	478 ± 7	487 ± 11
door-unlock	3421 ± 974	3277 ± 1265	2806 ± 1437	787 ± 1001	691 ± 356

Overall, we hope that this work will pave the way for a new direction in reinforcement learning safety research, shifting the focus from reward to agent behavior.

Thank you!