

# Summary of Three Recent Papers: Deep Reinforcement Learning and Adversarial Attacks

Presenter: Ji Gao



Department of Computer Science, University of Virginia

<https://qdata.github.io/deep2Read/>

# Paper List

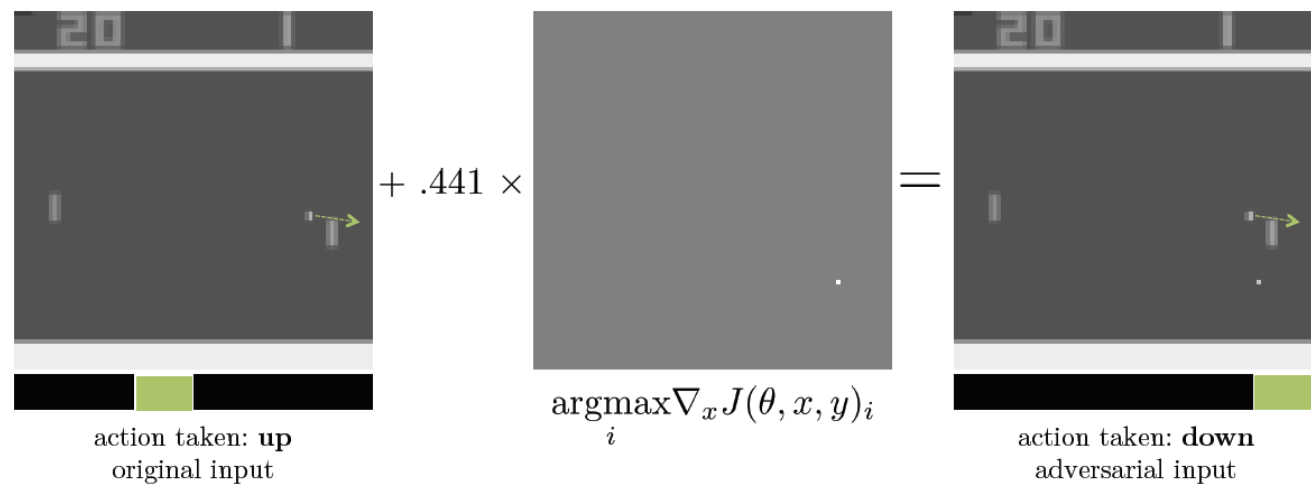
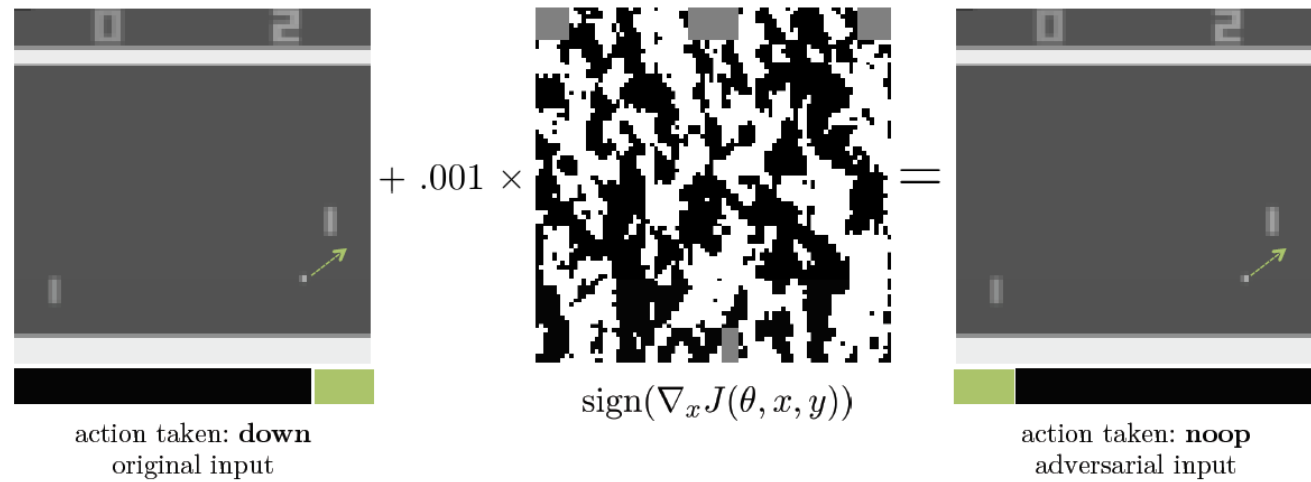
- Adversarial Attacks on Neural Network Policies
- Vulnerability of Deep Reinforcement Learning to Policy Induction Attacks
- Online Robust Policy Learning in the Presence of Unknown Adversaries
- Robust Deep Reinforcement Learning with Adversarial Attacks

# Adversarial Attacks on Neural Network Policies

*Sandy Huangy, Nicolas Papernotz, Ian Goodfellowx, Yan Duanyx, Pieter Abbeely*

- 2017 ICLR Workshop
- Idea: Craft adversarial samples in the input feature (State) of RL algorithm, lead to a large degrade of test-time performance
- Test-time

# Idea



# RL algorithm

- DQN: Deep Q Network
- TRPO: Trust Region Policy Optimization. Use a whole trajectory rollout on stochastic policy, penalized by the KL divergence between old and new policy.
- A3C: Asynchronous Advantage Actor-Critic. Asynchronous gradient descent on stochastic policy.

# Apply FGSM on Policy

- FGSM directly use  $J(\theta, x, y)$  to generate perturbation.
- In this case, just assume the policy  $\pi_\theta$  generated is good, and our goal is to flip the predicted policy.
- Note: Use an extra softmax on the DQN

# Different norms

$$\eta = \begin{cases} \epsilon \operatorname{sign}(\nabla_x J(\theta, x, y)) & \text{for constraint } \|\eta\|_\infty \leq \epsilon \\ \epsilon \sqrt{d} * \frac{\nabla_x J(\theta, x, y)}{\|\nabla_x J(\theta, x, y)\|_2} & \text{for constraint } \|\eta\|_2 \leq \|\epsilon \mathbf{1}_d\|_2 \\ \text{maximally perturb highest-impact dimensions with budget } \epsilon d & \\ & \text{for constraint } \|\eta\|_1 \leq \|\epsilon \mathbf{1}_d\|_1 \end{cases}$$

# Result

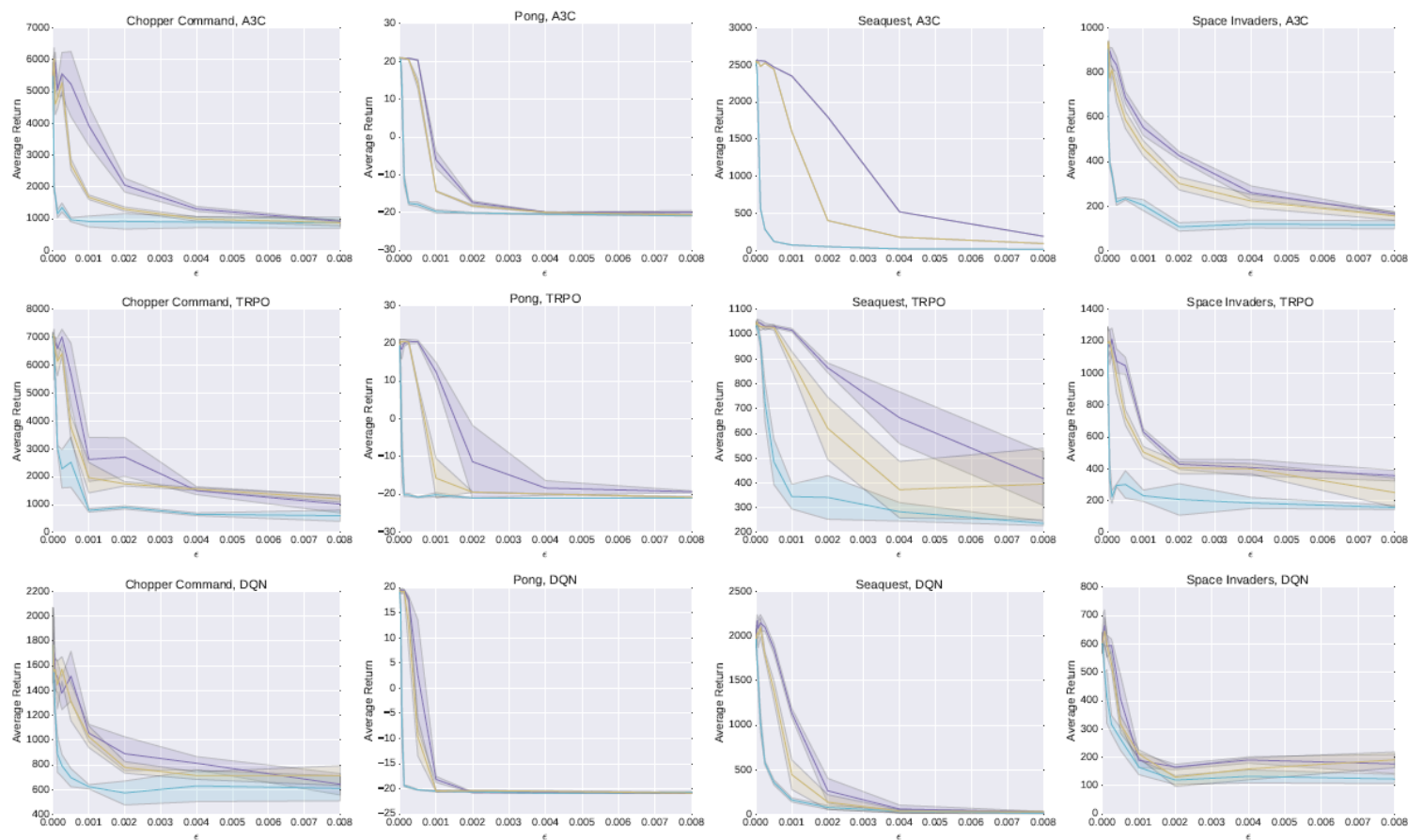


Figure 2: Comparison of the effectiveness of  $l_\infty$ ,  $l_2$ , and  $l_1$  FGSM adversaries on four Atari games trained with three learning algorithms. The average return is taken across ten trajectories. Constraint on FGSM perturbation:  $l_\infty$ -norm  $l_2$ -norm  $l_1$ -norm



# Transferability

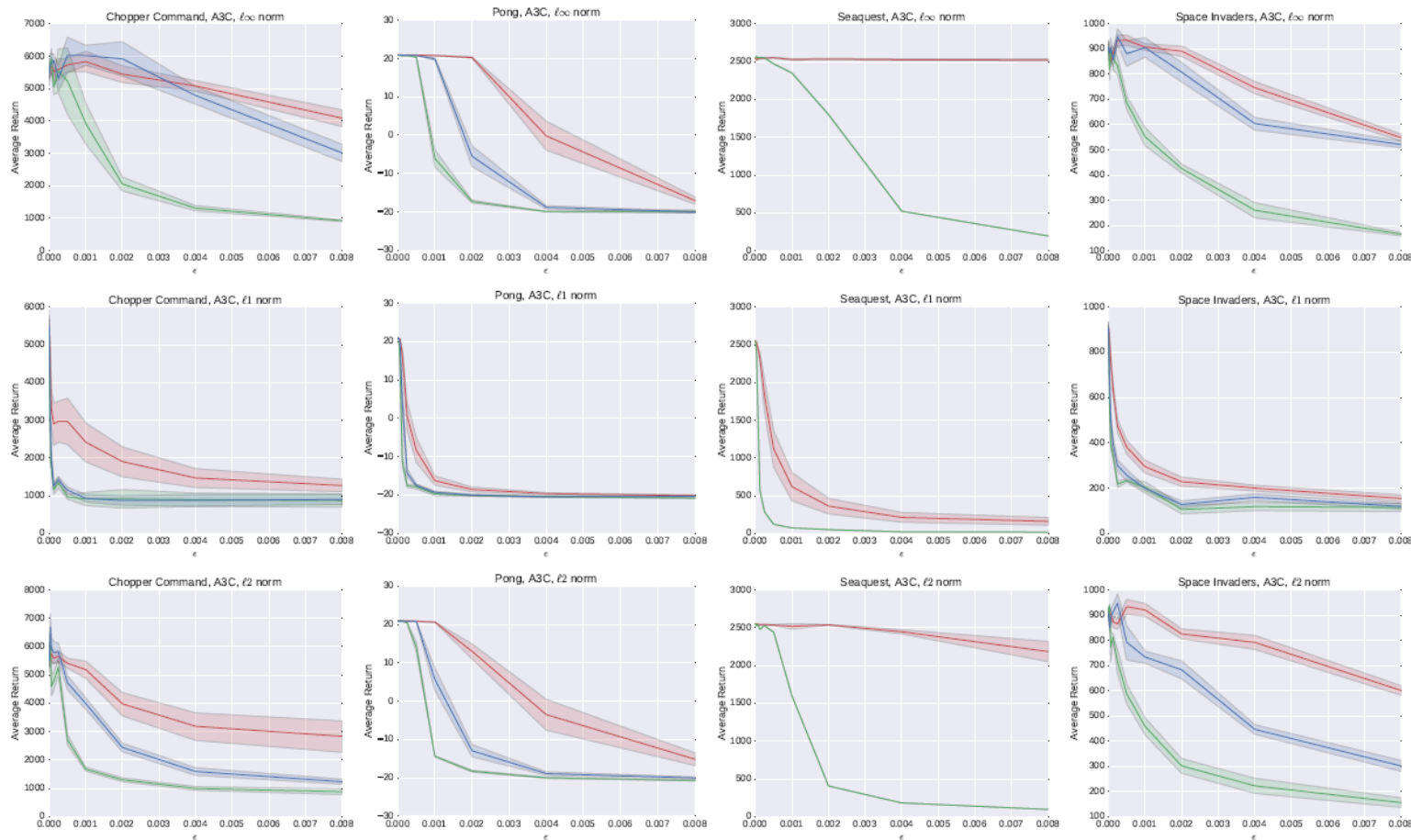


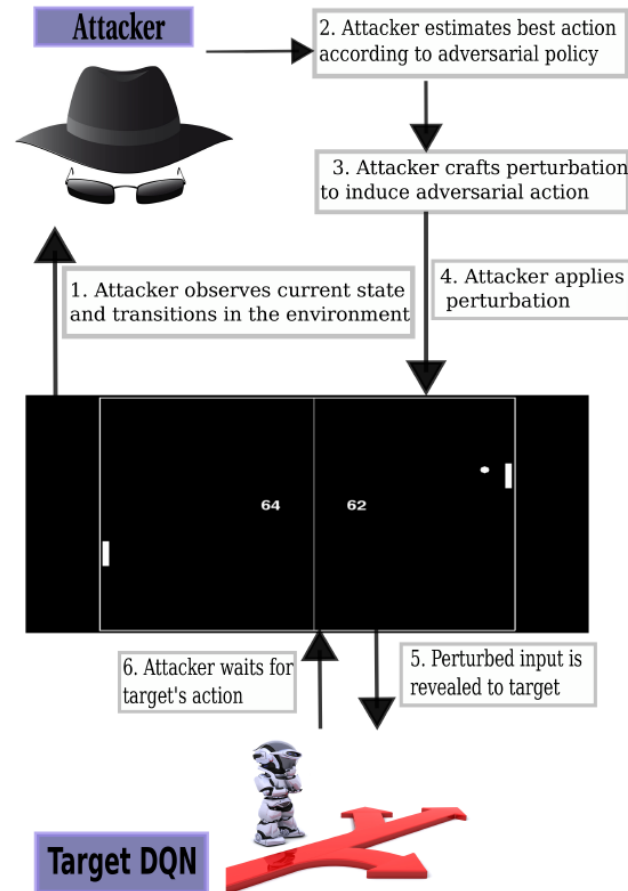
Figure 3: Transferability of adversarial inputs for policies trained with A3C. Type of transfer: ■ algorithm ■ policy ■ none

# Vulnerability of Deep Reinforcement Learning to Policy Induction Attacks

Vahid Behzadan and Arslan Munir

International Conference on Machine Learning and Data Mining in Pattern Recognition 2017

- Attack DQN



# Threat model

- Priors info: Structure of DQN, reward function  $R$
- Attacker: Work on state, not actions
- Magnitude of perturbation is smaller than  $\epsilon$

# Method

- Use another DQN to simulate DQN
- Use FGSM and JSMA

# Procedure

---

## Algorithm 1: Exploitation Procedure

---

**input** : adversarial policy  $\pi_{adv}^*$ , initialized replica DQNs  $Q'$ ,  $\hat{Q}'$ , synchronization frequency  $c$ , number of iterations  $N$

- 1 **for** *observation* = 1,  $N$  **do**
- 2     Observe current state  $s_t$ , action  $a_t$ , reward  $r_t$ , and resulting state  $s_{t+1}$
- 3     **if**  $s_{t+1}$  is not terminal **then**
- 4         set  $a'_{adv} = \pi_{adv}^*(s_{t+1})$
- 5         Calculate perturbation vector  $\hat{\delta}_{t+1} = Craft(\hat{Q}', a'_{adv}, s_{t+1})$
- 6         Update  $s_{t+1} \leftarrow s_{t+1} + \hat{\delta}_{t+1}$
- 7         Set  $y_t = (r_t + \max_{a'} \hat{Q}'(s_{t+1} + \hat{\delta}_{t+1}, a'; \theta'_-))$
- 8         Perform SGD on  $(y_t - Q'(s_t, a_t, \theta'))^2$  w.r.t  $\theta'$
- 9     **end**
- 10     Reveal  $s_{t+1}$  to target
- 11     **if** *observation* mod  $c = 0$  **then**  $\theta'_- \leftarrow \theta'$
- 12 **end**

---

# Experiment result

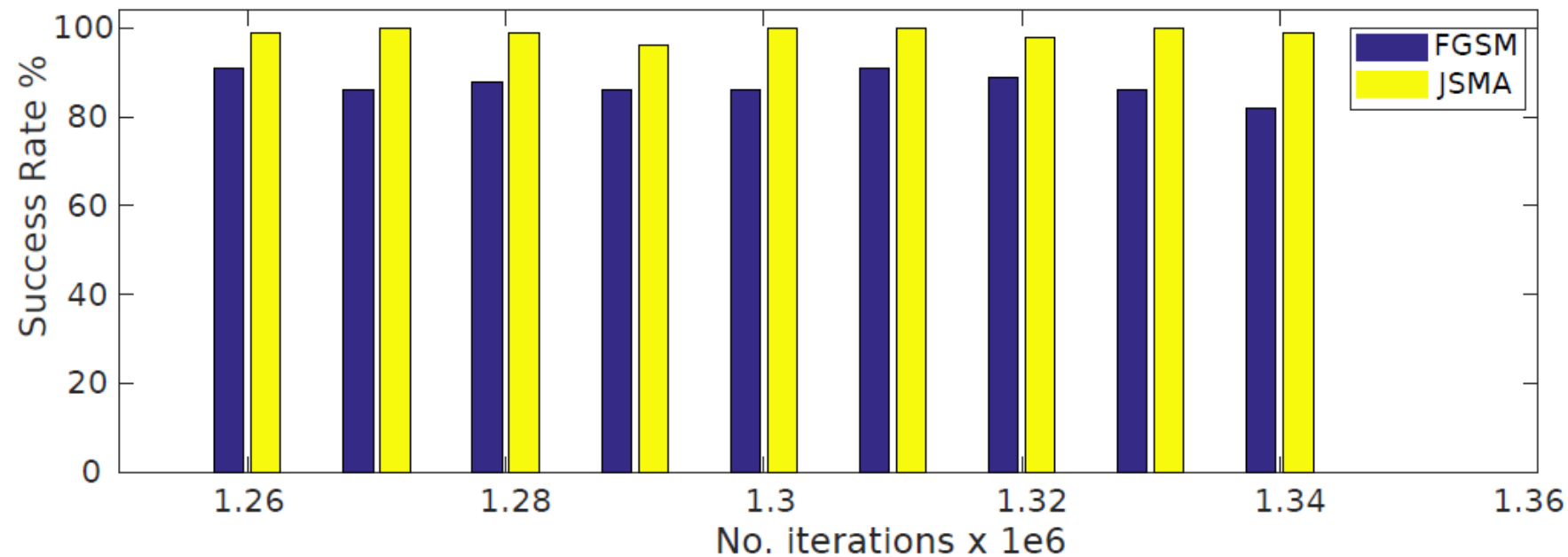
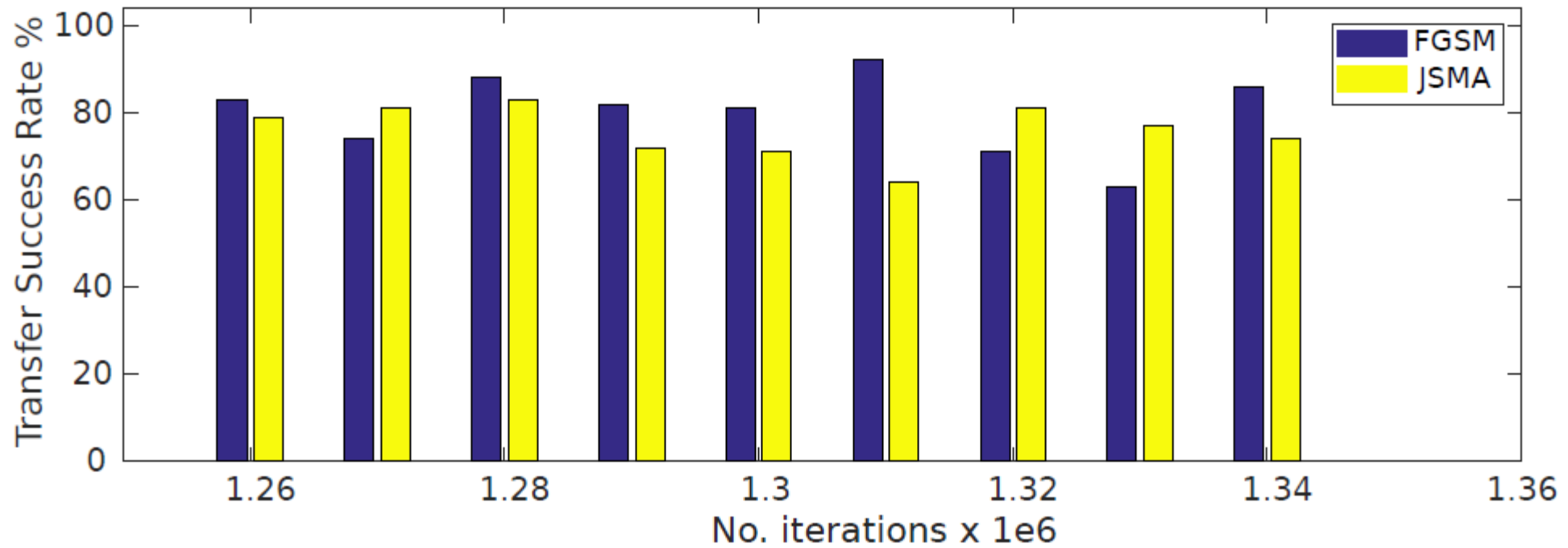


Fig. 4: Success rate of crafting adversarial examples for DQN

# Transferability



# Online Robust Policy Learning in the Presence of Unknown Adversaries

Aaron J. Havens, Zhanhong Jiang, Soumik Sarkar

- Setting: Finite-horizon discounted MDP
- Base Algorithm: TRPO
- Goal: Online mitigation of the adversarial perturbation

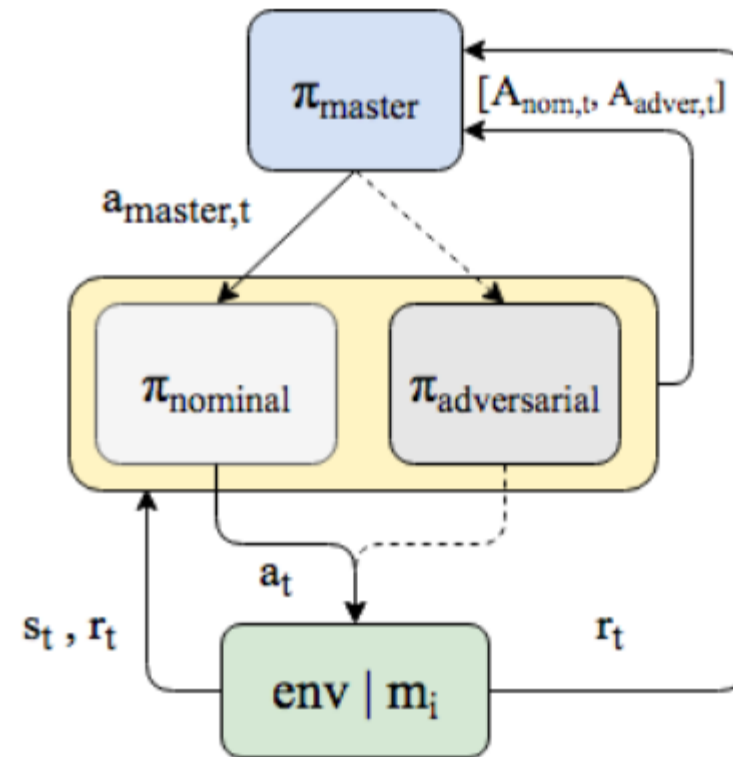


# Motivation

- In the learning, if the input changed, the expected reward will be different to the actual reward.
- In the online setting, an advantage (expect result – observed return) can be used as a good indicator that whether the input has been perturbed

# Method(General idea)

- Create multiple(2) policies
- A master policy can choose one policy from time to time.
- When it detects a change on the input, it will choose a different policy



(b) MLAH framework

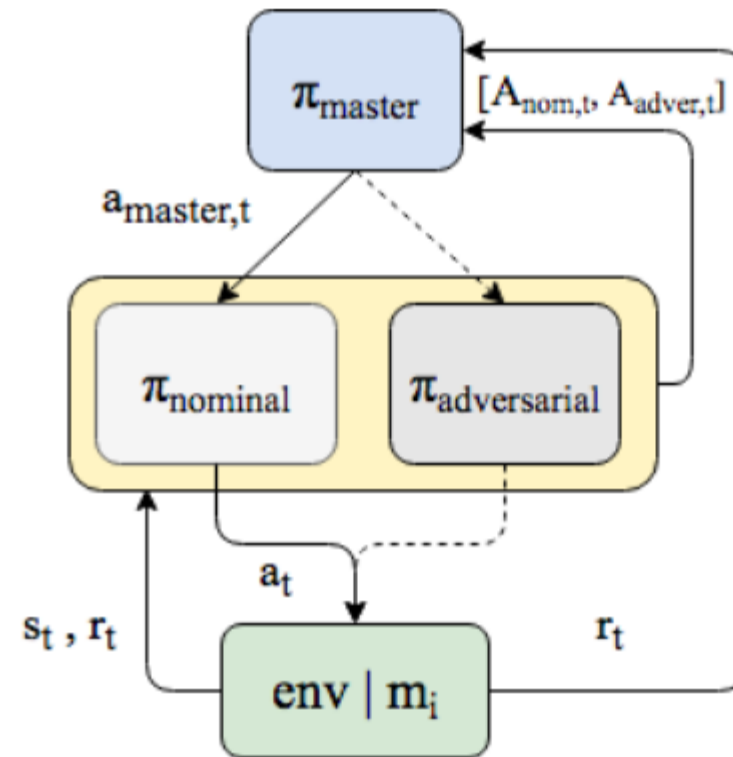
# Method

- Start with random policies
- For a step, estimate the advantage of both policies

$$\mathbf{A}_t = [A_{GAE,t-h}|\pi_{nom}, A_{GAE,t-h}|\pi_{adv}] \in \mathbb{R}^2$$

$$a_{master,t} = \pi_{*,t} = \operatorname{argmax}_a \mathbb{E}_{s_t, \pi_i, m_i \dots} \left[ \sum_{t=0}^T \gamma^t r(s_t, a) | m_i \right] \in \{\pi_{nom}, \pi_{adv}\}$$

- Optimize master policy according to the loss
- Optimize both policies



(b) MLAH framework

# Method

---

**Algorithm 1:** MLAH

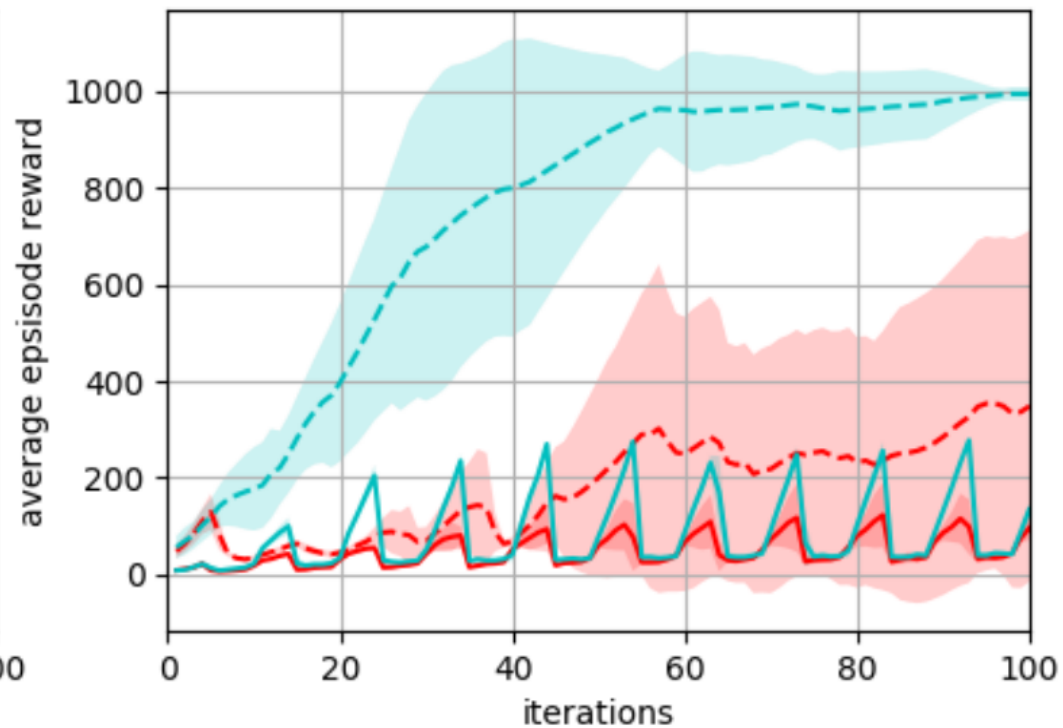
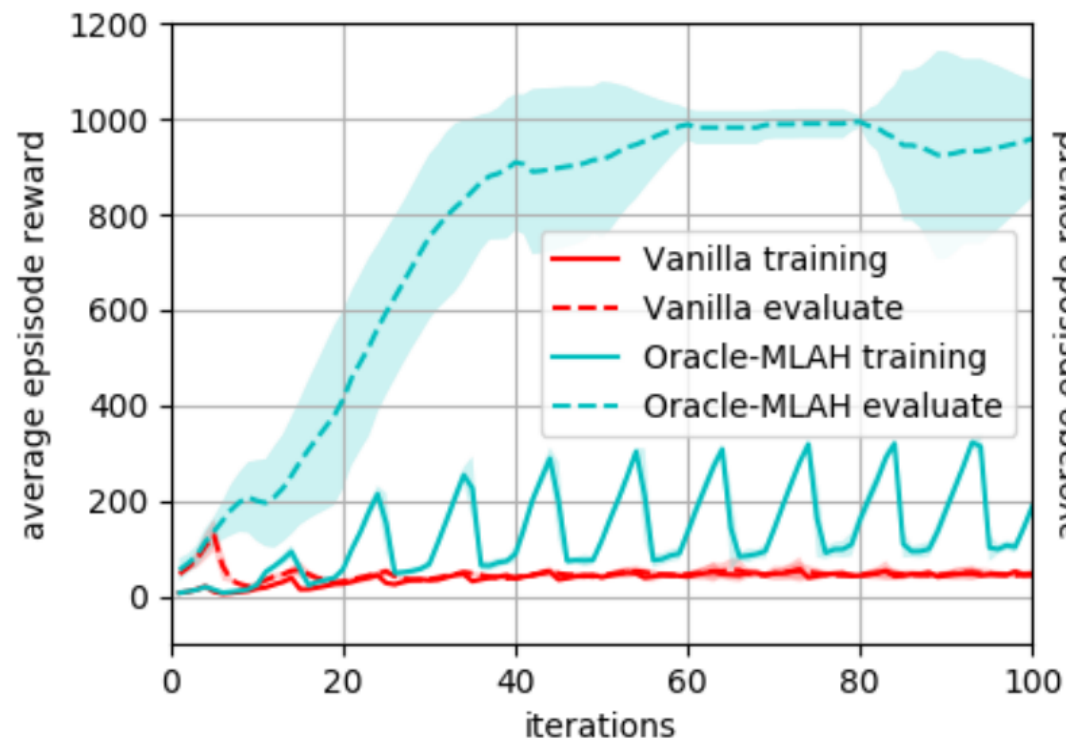
---

**Input:**  $\pi_{nom}$  and  $\pi_{adv}$  sub-policies parameterized by  $\theta_{nom}$  and  $\theta_{adv}$ ; Master policy  $\pi_{master}$  with parameter vector  $\phi$ .

- 1 Initialize  $\theta_{nom}, \theta_{adv}, \phi$
- 2 **for** *pre-training iterations [optional]* **do**
- 3 | Train  $\pi_{nom}$  and  $\theta_{nom}$  on only nominal experiences.
- 4 **end**
- 5 **for** *learning life-time* **do**
- 6 | **for** *Time steps  $t$  to  $t + T$*  **do**
- 7 | | Compute  $\mathbf{A}_t$  over sub-policies (see eq. 4)
- 8 | | select sub-policy to take action with  $\pi_{master}$  using  $\mathbf{A}_t$  as observations
- 9 | **end**
- 10 | Estimate all  $A_{GAE}$  for  $\pi_{nom}, \pi_{adv}$  over  $T$
- 11 | Estimate all  $A_{GAE}$  for  $\pi_{master}$  over  $T$  with respect to  $\mathbf{A}_t$  observations
- 12 | Optimize  $\theta_{nom}$  based on experiences collected from  $\pi_{nom}$
- 13 | Optimize  $\theta_{adv}$  based on experiences collected from  $\pi_{adv}$
- 14 | Optimize  $\phi$  based on all experiences with respect to  $\mathbf{A}_t$  observations
- 15 **end**

---

# Result



# Robust Deep Reinforcement Learning with Adversarial Attacks

Pattanaik, Anay, Zhenyi Tang, Shuijing Liu, Gautham Bommanan, and Girish Chowdhary.  
In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*

- Even naïve attack can degrade DRL algorithms
- Adversarial Training leads to significant increase in robustness to parameter variations for RL benchmarks
- Target RL algorithm:
  - DDQN
  - DDPG

# Perturbation

- Naïve attack: Random perturbation
- Gradient attack: Gradient based perturbation(FGSM)

# Adversarial attack

- Use beta distribution to sample noise from some distribution
- Sample multiple times and pick the largest one

---

**Algorithm 2** Naive attack (DDPG)

---

```
1: procedure NAIVE( $Q^{target}, U, s, \epsilon, n, \alpha, \beta$ )      ▷ Naive attack function takes trained target
   critic network  $Q^{target}$ , trained actor network  $U$ , current state(s), adversarial attack magnitude
   constraint( $\epsilon$ ), parameters of beta distribution( $\alpha, \beta$ ) and number of times to sample noise( $n$ ) as
   input
2:    $a^* = U(s), Q^* = Q^{target}(s, a^*)$       ▷ Determine optimal action and action value function
3:   for  $i = 1 : n$  do                                ▷ Sample a few times
4:      $n_i \sim \text{beta}(\alpha, \beta) - 0.5$           ▷ Sample noise
5:      $s_i = s + \epsilon * n_i$                     ▷ Possible adversarial state determined by sampled noise
6:      $a_{adv} = U(s_i)$                             ▷ Determine optimal action in potential adversarial state
7:      $Q_{adv}^{target} = Q^{target}(s, a_{adv})$       ▷ Determine the value of potential adversarial action
   corresponding to potential adversarial state for current state
8:     if  $Q_{adv}^{target} < Q^*$  then                ▷ if the potential adversarial state leads to bad action
9:        $Q^* = Q_{adv}^{target}$                     ▷ Store the value function of that potential bad action
10:       $s_{adv} = s_i$                             ▷ Store possible adversarial state
11:     else
12:       do nothing
13:     end if
14:   end for
15:   return  $s_{adv}$                                 ▷ Adversarial state
16: end procedure
```

---



# Adversarial training

---

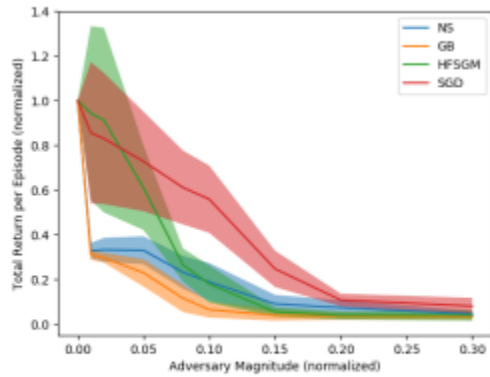
**Algorithm 5** Training with adversarial perturbation (DDQN)

---

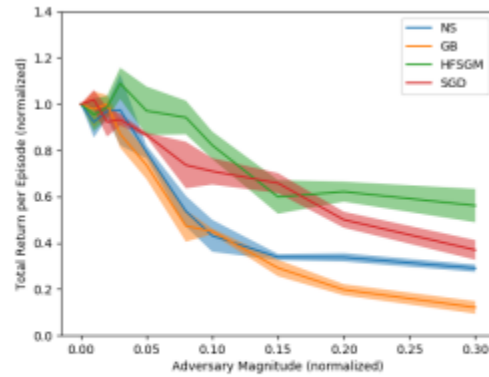
```
1: procedure ADV TRAIN ( $Q^{target}, Q$ )           ▷ Gradient based adversarial training method takes
   pre-trained network
2:   for  $i = 1 : iterations$  do                 ▷ Train adversarially for number of timesteps
3:     Reset the environment and receive observation
4:     while not terminal or not max time steps per episode reached do
5:        $s_{adv} = Grad(Q^{target}, Q, s, \epsilon, n, \alpha, \beta)$            ▷ Fool the agent
6:        $a = arg \max_a Q(s_{adv}, a)$        ▷ Fooled agent takes action according to behavior policy
7:        $s, r = Env(a, s)$  ▷ Environment returns next state and reward corresponding to state
    $s$  and action  $a$ 
8:       Update the weights of network according to DDQN algorithm
9:     end while
10:  end for
11: end procedure
```

---

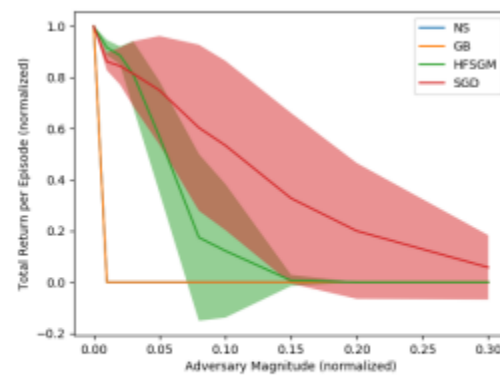
# Adversarial Attack performance



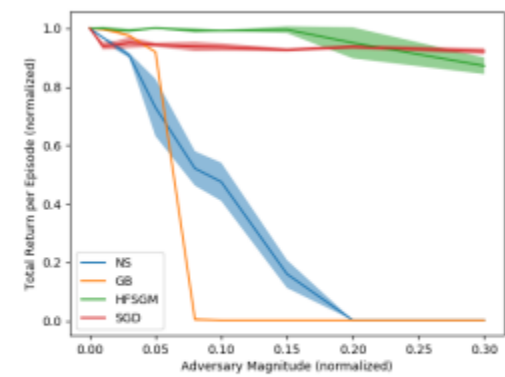
(a) DDQN Cart Pole



(b) RBF Q Cart Pole



(c) DDQN Mountain Car



(d) RBF Q mountain car

# Robust

- Test the adversarial trained agent on a wide range of parameters and compared it to “vanilla” DRL