

# Adversarial Policy Learning in Two-player Competitive Games

## S1 Proof of Proposition 1

**Proposition 1.** *In a two-player Markov game, if one agent follows a fixed policy, the state transition of the game system will depend only upon the policy of the other agent rather than their joint policies.*

**Proof.** Without the loss of generalizability, we denote the agent as  $\alpha$  and  $v$ , and assume the agent  $v$  follows a fixed policy. At state  $s_t$ , the probability of taking the joint actions  $(a_t^\alpha, a_t^v)$  and transiting to  $s_{t+1}$  is

$$\begin{aligned} P(s_{t+1}, a_t^\alpha, a_t^v | s_t) &= P(s_{t+1} | a_t^\alpha, a_t^v, s_t) P(a_t^\alpha, a_t^v | s_t) = P(s_{t+1} | a_t^\alpha, a_t^v, s_t) P(a_t^\alpha | a_t^v, s_t) \pi^v(a_t^v | s_t) \\ &= c \cdot P(s_{t+1} | a_t^\alpha, a_t^v, s_t) P(a_t^\alpha | a_t^v, s_t) = c \cdot P(s_{t+1} | a_t^\alpha, a_t^v, s_t) \pi^\alpha(a_t^\alpha | s_t), \end{aligned} \quad (1)$$

where  $P(a_t^v | s_t) = c$ . Given that at a time step  $t$ , an action of the agent  $a_t^\alpha$  depends only upon the current state  $s_t$ , we have  $\pi^\alpha(a_t^\alpha | s_t) = P(a_t^\alpha | a_t^v, s_t)$ .

As we can observe from Eqn. (1), during the adversarial training process, the only changed part is policy  $\pi^\alpha$ . As such, the change in  $\pi^\alpha$  determines the change in state transition and the in change both agent's value functions. To be specific, given a set of trajectories  $\{\tau_1, \dots, \tau_M\}$ , the state-value function of  $\alpha$  agent  $V_\pi^\alpha$  can be computed by

$$V_\pi^\alpha = \sum_{m=1}^M R^\alpha(\tau_m) P(\tau_m; \theta). \quad (2)$$

The state-value function of  $v$  agent  $V_\pi^v$  can be computed by

$$V_\pi^v = \sum_{m=1}^M R^v(\tau_m) P(\tau_m; \theta). \quad (3)$$

In Eqn. (2) and Eqn. (3),

$$\begin{aligned} P(\tau; \theta) &= P(s_0) \prod_{t=1}^T P(s_t, a_{t-1}^\alpha, a_{t-1}^v | s_{t-1}) \\ &= P(s_0) \prod_{t=1}^T P(s_t, a_{t-1}^\alpha | a_{t-1}^v, s_{t-1}) P(a_{t-1}^v | s_{t-1}). \end{aligned} \quad (4)$$

Since the victim agent follows a fixed policy,  $P(a_{t-1}^v | s_{t-1}) = c$ . Then, Eqn. (4) can be rewrote as

$$\begin{aligned} P(\tau; \theta) &= P(s_0) \prod_{t=1}^T P(s_t, a_{t-1}^\alpha, a_{t-1}^v | s_{t-1}) \\ &= P(s_0) \prod_{t=1}^T c \cdot P(s_t | a_{t-1}^\alpha, a_{t-1}^v, s_{t-1}) \pi^\alpha(a_{t-1}^\alpha | s_{t-1}). \end{aligned} \quad (5)$$

Similar to Eqn. (1),  $\pi^\alpha$  is the only changed component in Eqn. (5). Plugging Eqn. (5) into either Eqn. (2) or Eqn. (3), we can find out the change in both agent's state-value functions depend only upon the policy  $\pi^\alpha$ . Combining these observations with the aforementioned observation in Eqn. (1), we can conclude that the change in  $\pi^\alpha$  determines the change in state transition as well as the change in both agent's value functions.  $\square$

## S2 Proof of Lemma 1

**Lemma 1.** *Given a old policy  $\pi^\alpha$  and a new policy  $\pi^{\alpha'}$  in a two-agent Markov game, the difference of the victim state-value function under each policy is as follows.*

$$V_{\pi^{\alpha'}}^v(s) - V_{\pi^\alpha}^v(s) = E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t A_{\pi^\alpha}^v(s_t, a_t^\alpha) \right]. \quad (6)$$

**Proof.** Recall that in Section 3.3, we state that the victim value function can be redefined as follows

$$Q_{\pi^\alpha}^v(s_t, a_t^\alpha) = R^v(s_t, a_t) + \gamma E_{s_{t+1}|s_t, a_t} [V_{\pi^\alpha}^v(s_{t+1})]. \quad (7)$$

In addition, the tower property of conditional expectations gives the following equation

$$E_{X,Y}[f(x,y)] = E_X E_{Y|X}[f(x,y)] = E_{X,Y} E_{Y|X}[f(x,y)], \quad (8)$$

where  $x$  and  $y$  are random variables. Based on Eqn. (8), we also have

$$E_{\tau \sim \pi^{\alpha'}} V_{\pi^\alpha}^v(s_{t+1}) = E_{\tau \sim \pi^{\alpha'}} [E_{s_{t+1}|s_t, a_t} [V_{\pi^\alpha}^v(s_{t+1})]]. \quad (9)$$

Then, we can compute the victim state-value function difference.

$$\begin{aligned} V_{\pi^{\alpha'}}^v(s) - V_{\pi^\alpha}^v(s) &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t R^v(s_t, a_t) \right] - V_{\pi^\alpha}^v(s) \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t [R^v(s_t, a_t) - V_{\pi^\alpha}^v(s) + V_{\pi^\alpha}^v(s)] \right] - V_{\pi^\alpha}^v(s) \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t R^v(s_t, a_t) - \sum_{t=0}^{\infty} \gamma^t V_{\pi^\alpha}^v(s) + \sum_{t=0}^{\infty} \gamma^t V_{\pi^\alpha}^v(s) \right] - V_{\pi^\alpha}^v(s) \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t [R^v(s_t, a_t) - V_{\pi^\alpha}^v(s)] + \sum_{t=0}^{\infty} \gamma^{t+1} V_{\pi^\alpha}^v(s_{t+1}) + V_{\pi^\alpha}^v(s) \right] - V_{\pi^\alpha}^v(s) \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t [R^v(s_t, a_t) - V_{\pi^\alpha}^v(s)] + \sum_{t=0}^{\infty} \gamma^{t+1} V_{\pi^\alpha}^v(s_{t+1}) \right] \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t [R^v(s_t, a_t) - V_{\pi^\alpha}^v(s) + \gamma V_{\pi^\alpha}^v(s_{t+1})] \right] \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t [R^v(s_t, a_t) - V_{\pi^\alpha}^v(s) + \gamma E_{s_{t+1}|s_t, a_t} [V_{\pi^\alpha}^v(s_{t+1})]] \right] \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t [R^v(s_t, a_t) + \gamma E_{s_{t+1}|s_t, a_t} [V_{\pi^\alpha}^v(s_{t+1})] - V_{\pi^\alpha}^v(s)] \right] \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t [Q_{\pi^\alpha}^v(s_t, a_t^\alpha) - V_{\pi^\alpha}^v(s)] \right] \\ &= E_{\tau \sim \pi^{\alpha'}} \left[ \sum_{t=0}^{\infty} \gamma^t A_{\pi^\alpha}^v(s_t, a_t^\alpha) \right], \end{aligned} \quad (10)$$

where  $a_t = (a_t^\alpha, F^{\pi^v}(s_t))$ .

## S3 Proof of Theorem 1

**Theorem 1.** *The difference between  $V_{\pi^{\alpha'}}^v(s)$  and  $L_{\pi^\alpha}^v(\pi^{\alpha'})$  is bounded by:*

$$\begin{aligned} V_{\pi^{\alpha'}}^v(s) &\leq L_{\pi^\alpha}^v(\pi^{\alpha'}) + C_2 \mathbb{KL}^{\max}(\pi^\alpha || \pi^{\alpha'}) = M_{\pi^\alpha}^v(\pi^{\alpha'}), \\ C_2 &= \frac{4 \max_{s, a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)| \gamma}{(1-\gamma)^2}. \end{aligned} \quad (11)$$

**Proof.**

$$\begin{aligned}
& V_{\pi^{\alpha'}}^v(s) - L_{\pi^\alpha}^v(\pi^{\alpha'}) \\
&= \sum_s P(s_t = s | \pi^{\alpha'}) \sum_a \pi^{\alpha'}(a^\alpha | s) \sum_t \gamma^t A_{\pi^\alpha}^v(s, a^\alpha) - \sum_s P(s_t = s | \pi^\alpha) \sum_a \pi^\alpha(a^\alpha | s) \sum_t \gamma^t A_{\pi^\alpha}^v(s, a^\alpha) \\
&= \mathbb{E}_{s_t \sim \pi^{\alpha'}} \left[ \sum_t \gamma^t \mathbb{E}_{a \sim \pi^{\alpha'}(\cdot | s)} [A_{\pi^\alpha}^v(s, a^\alpha)] \right] - \mathbb{E}_{s_t \sim \pi^\alpha} \left[ \sum_t \gamma^t \mathbb{E}_{a \sim \pi^{\alpha'}(\cdot | s)} [A_{\pi^\alpha}^v(s, a^\alpha)] \right] \\
&= \sum_t \gamma^t [\mathbb{E}_{s_t \sim \pi^{\alpha'}} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] - \mathbb{E}_{s_t \sim \pi^\alpha} [\bar{A}_{\pi^{\alpha'}}^v(s_t)]],
\end{aligned} \tag{12}$$

where  $\bar{A}_{\pi^{\alpha'}}^v(s) = \mathbb{E}_{a \sim \pi^{\alpha'}(\cdot | s)} [A_{\pi^\alpha}^v(s, a^\alpha)]$ .

Let  $n_t$  denotes the number of time steps that  $a_i^{\alpha'} \neq a_i^\alpha$  for time step  $i < t$ , where  $a_i^{\alpha'} \sim \pi^{\alpha'}$  and  $a_i^\alpha \sim \pi^\alpha$ . That is, the number of time steps that  $\pi^{\alpha'}$  and  $\pi^\alpha$  disagrees before time step  $t$ .

Then,

$$\mathbb{E}_{s_t \sim \pi^{\alpha'}} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] - \mathbb{E}_{s_t \sim \pi^\alpha} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] = P(n_t > 0) (\mathbb{E}_{s_t \sim \pi^{\alpha'} | n_t > 0} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] - \mathbb{E}_{s_t \sim \pi^\alpha | n_t > 0} [\bar{A}_{\pi^{\alpha'}}^v(s_t)]). \tag{13}$$

Given that  $(\pi^{\alpha'}, \pi^\alpha)$  is an  $\beta$ -coupled policy pair [11], we have

$$P(a_i^{\alpha'} = a_i^\alpha) \geq 1 - \beta. \tag{14}$$

We change the original notation  $\alpha$  to  $\beta$  to avoid confusion with the  $\alpha$  defined in our paper (*i.e.*, the adversarial agent). Then, we have

$$p(n_t = 0) = \prod_{i=1}^t P(a_i^{\alpha'} = a_i^\alpha) \geq (1 - \beta)^t, \tag{15}$$

and  $p(n_t > 0) \leq 1 - (1 - \beta)^t$ . Then, we can derive

$$\begin{aligned}
& \mathbb{E}_{s_t \sim \pi^{\alpha'}} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] - \mathbb{E}_{s_t \sim \pi^\alpha} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] \\
&= P(n_t > 0) (\mathbb{E}_{s_t \sim \pi^{\alpha'} | n_t > 0} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] - \mathbb{E}_{s_t \sim \pi^\alpha | n_t > 0} [\bar{A}_{\pi^{\alpha'}}^v(s_t)]) \\
&\leq P(n_t > 0) (|\mathbb{E}_{s_t \sim \pi^{\alpha'} | n_t > 0} [\bar{A}_{\pi^{\alpha'}}^v(s_t)]| + |\mathbb{E}_{s_t \sim \pi^\alpha | n_t > 0} [\bar{A}_{\pi^{\alpha'}}^v(s_t)]|) \\
&\stackrel{(a)}{\leq} P(n_t > 0) 4\beta \max_{s, a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)| \\
&\leq (1 - (1 - \beta)^t) 4\beta \max_{s, a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)|.
\end{aligned} \tag{16}$$

Where (a) can be obtained based on the following relationship. First, given that  $\mathbb{E}_{a^\alpha \sim \pi^\alpha} [A_{\pi^\alpha}^v(s_t, a^\alpha)] = 0$ , we have

$$\begin{aligned}
\bar{A}_{\pi^{\alpha'}}^v(s_t) &= \mathbb{E}_{a^{\alpha'} \sim \pi^{\alpha'}(\cdot | s_t)} [A_{\pi^\alpha}^v(s_t, a^{\alpha'})] \\
&= P(a^\alpha \neq a^{\alpha'}) \mathbb{E}_{(a^\alpha, a^{\alpha'}) \sim (\pi^\alpha, \pi^{\alpha'}) | a^\alpha \neq a^{\alpha'}} [A_{\pi^\alpha}^v(s_t, a^{\alpha'}) - A_{\pi^\alpha}^v(s_t, a^\alpha)].
\end{aligned} \tag{17}$$

Based on Eqn. (17), we can further derive that

$$|\bar{A}_{\pi^{\alpha'}}^v(s_t)| \leq P(a^\alpha \neq a^{\alpha'}) \mathbb{E}[|A_{\pi^\alpha}^v(s_t, a^{\alpha'})| + |A_{\pi^\alpha}^v(s_t, a^\alpha)|] \leq \beta \cdot 2 \max_{s, a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)|. \tag{18}$$

Given that  $\bar{A}_{\pi^{\alpha'}}^v$  at any state fulfills the inequality in Eqn. (18), the expectation of  $\bar{A}_{\pi^{\alpha'}}^v$  over all the states also obeys this inequality. As such, we have

$$|\mathbb{E}_{s_t \sim \pi^\alpha | n_t > 0} [\bar{A}_{\pi^{\alpha'}}^v(s_t)]| \leq \beta \cdot 2 \max_{s, a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)|. \tag{19}$$

According to Eqn. (19), we can have the (a) in Eqn. (16).

Plugging Eqn. (16) into Eqn. (12), we have

$$\begin{aligned}
V_{\pi^{\alpha'}}^v(s) - L_{\pi^\alpha}^v(\pi^{\alpha'}) &= \sum_t \gamma^t [\mathbb{E}_{s_t \sim \pi^{\alpha'}} [\bar{A}_{\pi^{\alpha'}}^v(s_t)] - \mathbb{E}_{s_t \sim \pi^\alpha} [\bar{A}_{\pi^{\alpha'}}^v(s_t)]] \\
&\leq \sum_t \gamma^t (1 - (1 - \beta)^t) 4\beta \max_{s, a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)| \\
&= 4\varepsilon\beta \sum_t \gamma^t (1 - (1 - \beta)^t) \\
&= \frac{4\varepsilon\gamma\beta^2}{(1 - \gamma)(1 - \gamma(1 - \beta))} \\
&\leq \frac{4\varepsilon\gamma\beta^2}{(1 - \gamma)^2},
\end{aligned} \tag{20}$$

where  $\varepsilon = \max_{s,a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)|$ . According to a Proposition in [5],  $\beta = \max_s \mathbb{T}\mathbb{V}(\pi^\alpha(\cdot|s) || \pi^{\alpha'}(\cdot|s))$ . Then, we can derive

$$V_{\pi^{\alpha'}}^v(s) \leq L_{\pi^\alpha}^v(\pi^{\alpha'}) + \frac{4\varepsilon\gamma\beta^2}{(1-\gamma)^2}. \quad (21)$$

According to [10],  $\mathbb{T}\mathbb{V}(p||q)^2 \leq \mathbb{K}\mathbb{L}(p||q)$ . Plugging this relationship into Eqn. (21), we have

$$V_{\pi^{\alpha'}}^v(s) \leq L_{\pi^\alpha}^v(\pi^{\alpha'}) + C_2 \mathbb{K}\mathbb{L}^{\max}(\pi^\alpha(\cdot|s) || \pi^{\alpha'}(\cdot|s)), \quad (22)$$

where  $C_2 = \frac{4\gamma \max_{s,a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)|}{(1-\gamma)^2}$ .  $\square$

## S4 Adversarial Learning Algorithm

In this section, we first describe how to transform the approximated objective function  $M_{\pi^\alpha}(\pi^{\alpha'})$  into our final adversarial learning objective function, followed by our adversarial learning algorithm.

Here, we first rewrite the Eqn. (11) in Section 3.3,

$$M_{\pi^\alpha}(\pi^{\alpha'}) = \sum_s \rho_{\pi^\alpha}(s) \sum_a \pi^{\alpha'}(a^\alpha|s) (A_{\pi^\alpha}^\alpha(s, a^\alpha) - A_{\pi^\alpha}^v(s, a^\alpha)) - C \mathbb{K}\mathbb{L}^{\max}(\pi^\alpha || \pi^{\alpha'}) + C_3, \quad (23)$$

where  $C = C_1 - C_2$  and  $C_3 = (V_{\pi^\alpha}^\alpha(s) - V_{\pi^\alpha}^v(s))$  are constants. Then, by following the method introduced in TRPO, we can further transform the maximization of Eqn. (23) into the following form

$$\begin{aligned} & \text{maximize}_{\pi^{\alpha'}} \sum_s \rho_{\pi^\alpha}(s) \sum_a \pi^{\alpha'}(a^\alpha|s) (A_{\pi^\alpha}^\alpha(s, a^\alpha) - A_{\pi^\alpha}^v(s, a^\alpha)), \\ & \text{s.t. } \mathbb{E}_{s \sim \pi^\alpha} [\mathbb{K}\mathbb{L}(\pi^\alpha(\cdot|s) || \pi^{\alpha'}(\cdot|s))] \leq \delta. \end{aligned} \quad (24)$$

As we can see from the equation above, this transformation replaces  $\mathbb{K}\mathbb{L}^{\max}(\pi^\alpha || \pi^{\alpha'})$  in Eqn. (23) with  $\mathbb{E}_{s \sim \pi^\alpha} [\mathbb{K}\mathbb{L}(\pi^\alpha || \pi^{\alpha'})]$  for the following reasons.  $\mathbb{E}_{s \sim \pi^\alpha} [\mathbb{K}\mathbb{L}(\pi^\alpha || \pi^{\alpha'})]$  is the average KL divergence between  $\pi^\alpha$  and  $\pi^{\alpha'}$ , which can be easily computed by using the Monte Carlo method [14]. Using this expectation as the substitution for maximum KL divergence, it is no longer required us to perform intensive computation at each state. In addition to the computation benefit, the replacement of maximum KL divergence indicates the ease of solving optimization. When performing optimization with maximum KL divergence, we have to introduce a constraint for each state. Given a two-player game with many states, this means imposing a large number of constraints on our optimization problem and potentially introduces the difficulty in getting an optimal solution. Note that, as is experimented in [9, 8], applying such an approximation imposes only a minor variation to the resolved policy.

As we can also observe from Eqn. (23), in Eqn. (24), we also transform the term  $\text{maximize} - C \mathbb{K}\mathbb{L}^{\max}(\pi^\alpha || \pi^{\alpha'})$  into a trust region constraint  $\mathbb{E}[\mathbb{K}\mathbb{L}(\pi^\alpha || \pi^{\alpha'})] \leq \delta$ .<sup>1</sup> As is discussed in [11], this transformation could enable a larger step size for the optimization process and thus accelerate the optimization process.

Even with all the transformation above, optimizing Eqn. (24) is still challenging. As we can see, this optimization objective involves the computation of  $\sum_a \pi^{\alpha'}(a^\alpha|s)$  which contains the actions tied to the new policy  $\pi^{\alpha'}$ . Before getting the optimization result, these actions are unknown. As a result, computing  $\sum_a \pi^{\alpha'}(a^\alpha|s) (A_{\pi^\alpha}^\alpha(s, a^\alpha) - A_{\pi^\alpha}^v(s, a^\alpha))$  is intractable.

To solve this problem, we again follow the idea of TRPO, apply an important sampling estimator

$$\begin{aligned} & \sum_a \pi^{\alpha'}(a^\alpha|s) (A_{\pi^\alpha}^\alpha(s, a^\alpha) - A_{\pi^\alpha}^v(s, a^\alpha)) \\ &= \sum_a \frac{\pi^{\alpha'}(a^\alpha|s)}{\pi^\alpha(a^\alpha|s)} \pi^\alpha(a^\alpha|s) (A_{\pi^\alpha}^\alpha(s, a^\alpha) - A_{\pi^\alpha}^v(s, a^\alpha)) \\ &= \mathbb{E}_{a \sim \pi^\alpha} \left[ \frac{\pi^{\alpha'}(a^\alpha|s)}{\pi^\alpha(a^\alpha|s)} (A_{\pi^\alpha}^\alpha(s, a^\alpha) - A_{\pi^\alpha}^v(s, a^\alpha)) \right], \end{aligned} \quad (25)$$

<sup>1</sup> $C = \frac{4\gamma(\max_{s,a^\alpha} |A_{\pi^\alpha}^\alpha(s, a^\alpha)| - \max_{s,a^\alpha} |A_{\pi^\alpha}^v(s, a^\alpha)|)}{(1-\gamma)^2}$ . Given that, in each iteration, our optimization maximizes the value function difference between the adversary and victim, we can obtain  $C > 0$  for the current policy  $\pi^\alpha$ .

---

**Algorithm 1:** Adversarial learning algorithm.

---

- 1 **Input:** The adversarial policy  $\pi_\theta^\alpha$  parameterized by  $\theta$ , the state-value function  $V_{\pi^\alpha}^\alpha$  and  $V_{\pi^\alpha}^v$ ,
  - 2 with parameters  $v_\alpha$  and  $v_v$ , respectively.
  - 3 **Initialization:** Initialize  $\theta^{(0)}$ ,  $v_\alpha^{(0)}$  and  $v_v^{(0)}$ .
  - 4 **for**  $k = 0, 1, 2, \dots, K$  **do**
  - 5     Use the current adversarial policy  $\pi_{\theta^{(k)}}^\alpha$  to play with the victim agent with a fixed policy  $\pi^v$ , and collect a set of trajectories  $\mathcal{D}^{(k)} = \{\tau_i\}$ , where  $i = 1, 2, \dots, |\mathcal{D}^{(k)}|$ .
  - 6     For each trajectory  $\tau_i$ , compute the advantage at each time step  $t$  ( $t = 0, 1, 2, \dots, |\tau_i|$ ):
  - 7      $A_{i_t}^{\alpha(k)} = r_{i_t}^{\alpha(k)} + \gamma V^{\alpha(k)}(o_{i_{t+1}}^{\alpha(k)}) - V^{\alpha(k)}(o_{i_t}^{\alpha(k)})$ ;
  - 8      $A_{i_t}^{v(k)} = r_{i_t}^{v(k)} + \gamma V^{v(k)}(o_{i_{t+1}}^{\alpha(k)}) - V^{v(k)}(o_{i_t}^{\alpha(k)})$ ,
  - 9     where we omit the subscript  $\pi_{\theta^{(k)}}^\alpha$  for simplicity.
  - 10    Introduce  $A_{i_{0:|\tau_i|}}^{\alpha(k)}$  and  $A_{i_{0:|\tau_i|}}^{v(k)}$  ( $i = 1 : |\mathcal{D}^{(k)}|$ ) into Eqn. (27) and obtain a new policy by maximizing the objective function in Eqn. (27).
  - 11    Update  $v_\alpha^{(k)}$  and  $v_v^{(k)}$  by solving  $\operatorname{argmin} \frac{1}{T} \sum_{t=0}^T (V(o_t) - (r_t + \gamma V_{old}(o_{t+1})))^2$ .
  - 12 **end**
  - 13 **Output:** The adversarial policy network  $\pi_\theta^\alpha$ .
- 

and thus transform Eqn. (24) into the form of

$$\begin{aligned} & \operatorname{argmax}_\theta \mathbb{E}_{\pi_{old}^\alpha} \left[ \frac{\pi_\theta^\alpha(a_t^\alpha | s_t)}{\pi_{old}^\alpha(a_t^\alpha | s_t)} (A_{\pi_{old}^\alpha}^\alpha(a_t^\alpha, s_t) - A_{\pi_{old}^\alpha}^v(a_t^\alpha, s_t)) \right], \\ & s.t. \mathbb{E}_{s_t \sim \pi_{old}^\alpha} [\mathbb{KL}(\pi_{old}^\alpha(\cdot | s_t) || \pi_\theta^\alpha(\cdot | s_t))] \leq \delta. \end{aligned} \quad (26)$$

As we can observe from the equation above, the expectation does not rely upon the actions pertaining to the new policy ( $\pi_\theta^\alpha$ ) but those tied to the old one ( $\pi_{old}^\alpha$ ). To learn an adversarial policy, we can optimize this objective function by using the algorithm introduced in TRPO [11]. However, in order to further improve the efficiency and effectiveness of the learning process, we follow the PPO algorithm [12], apply the clipped ratio operation, and obtain the following optimization function

$$\begin{aligned} & \operatorname{argmax}_\theta \mathbb{E}_{(a_t^\alpha, s_t) \sim \pi_{old}^\alpha} [\min(\operatorname{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) A_t^\alpha, \rho_t A_t^\alpha) - \min(\operatorname{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) A_t^v, \rho_t A_t^v)], \\ & \rho_t = \frac{\pi_\theta^\alpha(a_t^\alpha | s_t)}{\pi_{old}^\alpha(a_t^\alpha | s_t)}, A_t^\alpha = A_{\pi_{old}^\alpha}^\alpha(a_t^\alpha, s_t), A_t^v = A_{\pi_{old}^\alpha}^v(a_t^\alpha, s_t). \end{aligned} \quad (27)$$

In this work, we use this equation as our ultimate objective function and follow the procedure below to resolve this objective. Algorithm 1 shows our proposed adversarial learning algorithm. Specifically, we first approximate the corresponding state-value functions by using two deep neural networks, at the time step  $t$ , each of which takes as input the adversarial observation  $o_t^\alpha$  and outputs the approximated value for the state-value function  $V_t^\alpha$  and  $V_t^v$ . Second, we compute the parameters of these two networks by solving the optimization function in line 11. With the parameters resolved, we further update the adversarial policy network by solving the optimization function in Eqn. (27). In each training iteration, we gather a set of training trajectories by using the current adversarial agent to play with the fixed victim agent. By using the collected trajectories, we update the adversarial policy network and the networks pertaining to the two state-value functions. Note that, compared with the PPO algorithm used in the state-of-art attack [4], our proposed learning algorithm trains one additional value function and solves a more complicated optimization function. This leads to extra computational cost. To estimate this extra cost, we use the same machine (a server with 32 CPUs) to run our method and the state-of-art attack and record their runtimes. The average runtime of our method is about 1.4X over that of the existing attack (*e.g.*, 20 hours vs. 16 hours on the You-Shall-Not-Pass game, 32 hours vs. 23 hours on the Kick-And-Defend game). Considering the significant improvement in exploitability, we believe that this amount of extra cost is acceptable.

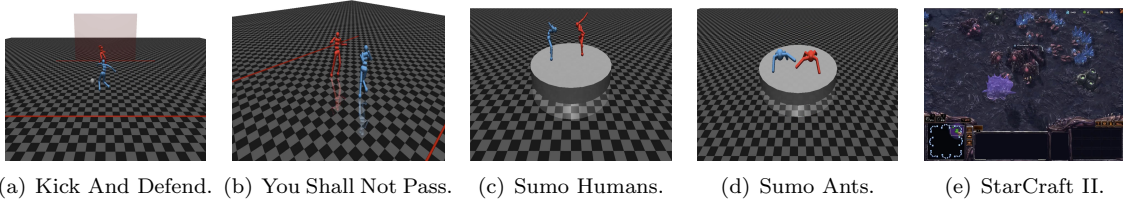


Figure S1: The snapshots of the two-player games selected for our experiment. The first four games are from MuJoCo game zoo and the last one is StarCraft II.

## S5 Implementation and Experiment Setups

### S5.1 Implementation and Hyper-parameter selection

We implemented our learning algorithm based on the TensorFlow [1] and Stable Baselines [7] packages. We implemented the baseline attack [4] based on the code released by the authors: <https://github.com/HumanCompatibleAI/adversarial-policies>. We also implemented the game environment wrappers based on the OpenAI Gym (*i.e.*, <https://gym.openai.com/>), Multi-agent Competition (*i.e.*, <https://github.com/openai/multiagent-competition>), as well as PySC2 Extension [13] packages.

In the following, we specify the choices of hyper-parameters for our attack and the state-of-art attack [4]. Both attacks have two sets of hyper-parameters: policy network/value function architectures and learning algorithm hyper-parameters (*e.g.*, clipping parameter  $\epsilon$ , discount factor  $\gamma$ , learning rate). To ensure a fair comparison, we adopted the same set of hyper-parameters in these two attacks. To be specific, for the adversarial policy network/state-value function architecture, we followed the choice in [4] and set them as Multilayer Perceptron with different layers for different games. The details of the network architectures can be found in <https://github.com/HumanCompatibleAI/adversarial-policies>. Regarding the learning algorithm hyper-parameters, we also used the default choice of the state-of-art attack [4]. The exact values are also shown in <https://github.com/HumanCompatibleAI/adversarial-policies>. Regarding the adversarial retraining experiments, we directly retrained the original victim agents with their original policy network/state-value function architectures and the same set of training hyper-parameters used in the attack experiments.

For the StarCraft II game, we also used the same set of hyper-parameters for these two attacks. Specifically, we directly adopted the default choices released by the PySC2 Extension platform, upon which we train the RL agents. The network architectures are also Multilayer Perceptrons. The detailed network architectures and the value of the training hyper-parameters can be found in <https://github.com/Tencent/TStarBot1>. It should be noted that compared to state-of-the-art attack, the only additional hyper-parameters introduced by our attack is the weight between the first and second term in our learning objective function. We varied the weight of the second term between [1, 4] and found that this variation imposes only a minor change upon the exploitability and transferability of our attack. As such, we gave these two terms the equal weight in our experiments.

### S5.2 Experiment Setups

**Game selection & obtaining victim agents.** We select five games for our experiments, including four robotic games from MuJoCo game zoo [15] and one real-time strategy game – StarCraft II [16]. Researchers commonly adopt these games in academia and industry to evaluate reinforcement learning algorithms in a two-player game context (*e.g.*, [2, 4, 6, 13]). For each of the games, researchers have released many benchmark game bots [2, 16]. Concerning the bots designed for MuJoCo games, they are all trained through DRL. However, the policy networks used in the bots are different (*e.g.*, “Sumo-Humans” and “You-Should-Not-Pass” use LSTM and an MLP as their policy networks, respectively). In this work, we use the following criteria to select our victim agent. First, we train an adversarial agent by using an existing technique [4]. Then, we use it to play with each of the agents and record the winning rate. For the agent demonstrating the highest winning rate against the adversary, we choose

it as the victim agent of that game. For the agent with the second-highest winning rate, we select it as a regular agent for evaluating the transferability of our adversarial agent and the generalizability of our retrained victim agent. It should be noted that we compute the winning rate by having the corresponding agent play with its opponent for 100 rounds and reporting the number of its wins. It should also be noted that for the asymmetric games “You-Shall-Not-Pass” and “Kick-And-Defend”, we select the runner and kicker as the victim agent, respectively.

Regarding the real-time strategy game StarCraft II, the game vendor, and their collaborators release seven bots indicating different master levels. (*i.e.*, level-1 to level-7 represents amateur to elite). These bots take actions under the guidance of different sets of pre-defined rules but not a policy network trained with an RL algorithm. As a result, we follow the method proposed by DeepMind [3], use the PPO algorithm to prepare two game agents, and ensure both of our game agents could demonstrate the decisive winning rates (*i.e.*, > 94%) against all the rule-based agents. In this work, we employ one agent as the victim agent and the other as the regular agent for the StarCraft game. In the following, we provide a more detailed description of each of the games mentioned above. Upon the acceptance of this work, we will release our source code and all the agents/environments used for our evaluation.

**MuJoCo–Kick And Defend.** This is a soccer penalty shootout, in which the kicker (*i.e.*, the blue humanoid robot in Figure 1(a)) intends to shoot the ball into the net (*i.e.*, the grey region on the red line in Figure 1(a)), whereas the defender (*i.e.*, the red humanoid robot in Figure 1(a)) prevents the kicker from scoring the goal. A successful scoring gives the kicker +1000 reward and the defender opponent -1000 reward. On the contrary, A successful defending gives the defender +1000 reward and the kicker -1000 reward. The defender is awarded an additional +500 reward if it saves a penalty and establishes certain desired behaviors. However, if the defender moves out of a defined goalkeeping region during a game, it gets a punishment of -1000 reward, and the game will end as a draw. Note that, in this game, a game episode exceeding the maximum time is treated as a successful defending.

**MuJoCo–You Shall Not Pass.** As is illustrated in Figure 1(b), the two agents in this game start by facing each other. Then, the blue humanoid robot (*i.e.*, runner) starts to run towards the finish line (*i.e.*, indicated by the red line in Figure 1(b)). Meanwhile, the red humanoid robot (*i.e.*, blocker) attempts to block the blue robot from reaching the finish line. If the red robot successfully stops its opponent and it keeps standing till the end of a game, it could receive +1000 reward. If it blocks its opponent but falls into the ground, it gets 0 reward. In both cases, the blue robot gets -1000 reward. On the contrary, if the blue robot reaches the finish line, it receives +1000 reward, and the red robot gets -1000 reward.

**MuJoCo–Sumo Humans and Sumo Ants.** In both games, the robots are randomly initialized at different positions on the grey round arena in Figure 1(c) and 1(d). Then, they start to move and push each other. One of the agents wins if it knocks its opponent into the ground or pushes it out of the arena. The winner receives +1000 reward, and the loser gets -1000 reward. If one agent falls from the arena without contacting its opponent or the game exceeds the maximum time, the game ends with a tie. Different from the games mentioned above, where the agent receives 0 reward in a tie game, in Sumo games, both agents get a penalty of -1000 reward for a draw. As is shown in Figure 1(c) and 1(d), the only difference between Sumo Humans and Sumo Ants is the shape of the robots. Note that, different from the two games introduced above, the agents are symmetric in the Sumo games.

**StarCraft II.** As is depicted in Figure 1(e), the base of each player is randomly placed at a corner on the map. Then, the players start to take action according to their strategies. The goal for each player is to defeat its opponent within a limited time. A game exceeding the time limit ends as a tie. In this paper, we train the reinforcement learning agents (players) on the `PySC2 Extension` platform released by [13]. To be specific, it designs 165 macro actions for an agent, each of which is a combination of the original operations in StarCraftII games. These macro actions can be categorized into five types – collecting resources, constructing buildings, producing workers and soldiers, upgrading technology, and combating. More details about the macro actions can be found in [13]. At the end of a game, each agent receives a reward based on the game result: 1 (win), 0 (tie), and -1 (lose). During the game, they also receive some additional rewards based on the number of enemies they have killed and the amount of resources they collected. Similar to [13], we consider a two-player competitive full-game in our experiments, in which both players belong to Zerg. Training an RL agent on a real gaming map

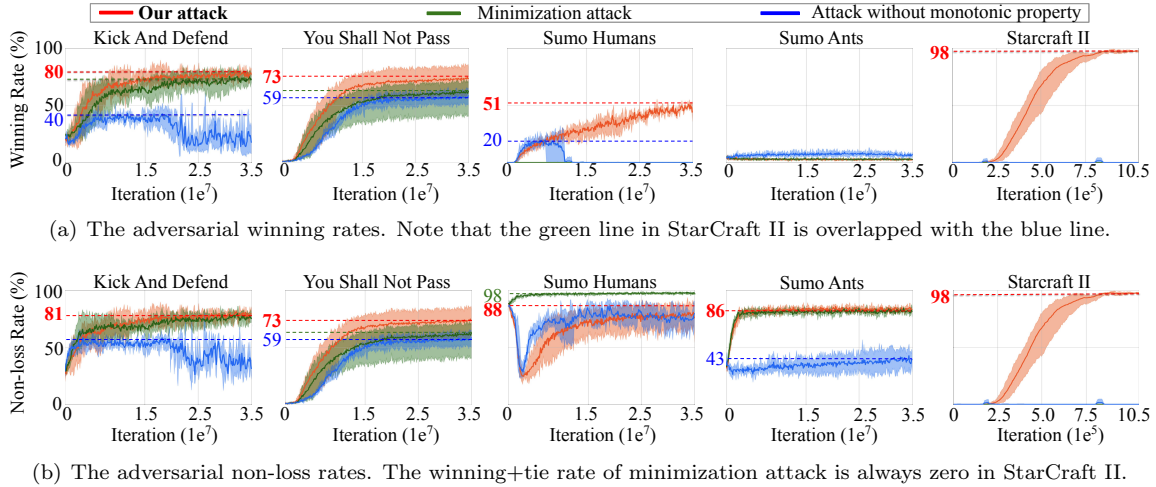


Figure S2: The performance comparison of adversarial agents trained with our attack and other two comparison methods: the attack that only minimizes the victim value function indicated by “Minimization attack” (*i.e.*, green lines and shadows in the figure) and the attack without monotonic property (*i.e.*, blue lines and shadows in the figure).

Table S1: Victim agents’ performances against our adversarial agents before retraining.

	Kick And Defend	You Shall Not Pass	Sumo Humans	Sumo Ants	StarCraft II
Winning (%)	25.0	30.0	30.0	15.0	2.0
Non-loss (%)	26.0	30.0	63.0	97.0	2.0

requires a large amount of time and computational resources. It takes [13] about two days and more than 3,000 CPUs to train an agent in a real gaming map. Due to limited computation resources, we use a smaller map designed specifically for training RL agents [16] rather than the real gaming maps. As is demonstrated in [6], the results of the smaller maps can be generalized to the real one by training the agents for a longer time.

**(Re)training & performance quantification.** From Algorithm 1, we could quickly note that, in each of the training iterations, our adversarial agent interacts with the victim, collects trajectories, and updates its policy network accordingly. To reduce the impact of state randomness (*e.g.*, the agent’s initial position on the map and the probabilistic state transitions) upon our adversarial agent’s performance, we follow the previous research [2, 4] and repeat each experiment six times by randomly selecting different initial states. With this setup, we report the average performance of our adversarial agents as well as their performance variance. Also, it should be noted that the algorithm updates our policy at each iteration. Therefore, we report the average performance of an adversarial agent every time its policy is updated. In our training process, we keep updating our adversarial agent iteratively until the agent performance converges. For all MuJoCo games, in the training process, the adversarial agent performance converges after 35 million iterations. For StarCraft II, it plateaus after 1.05 million iterations. As we will discuss below, we also design an experiment to evaluate the effectiveness of adversarial retraining for victim agents. For MuJoCo and StarCraft games, in the process of victim agent retraining, agent performance stays stable after 10 million and 2.2 million iterations, respectively. Similar to the setup for learning an adversarial agent, when retaining a victim agent, we follow the same setup process, which repeats our experiment for six times and reports average performance accordingly.

## S6 Additional Experiments.

**Significance of performance difference.** Section 4.1 visualizes the winning rate comparison between our adversarial agents and that obtained by the baseline approach. Here, we also conduct



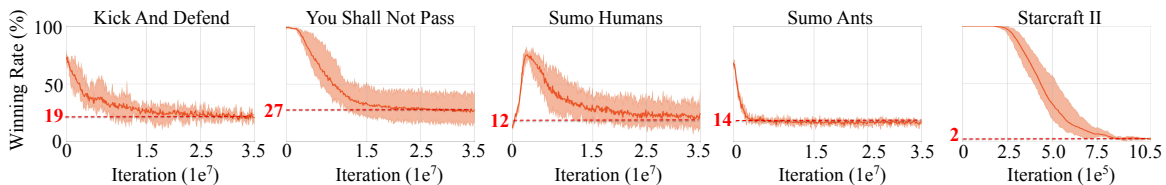


Figure S3: The changes in the winning rates of victim agents when training our adversarial agents against them.

Table S2: Mean, std and the  $p$ -value of the winning rate diff.

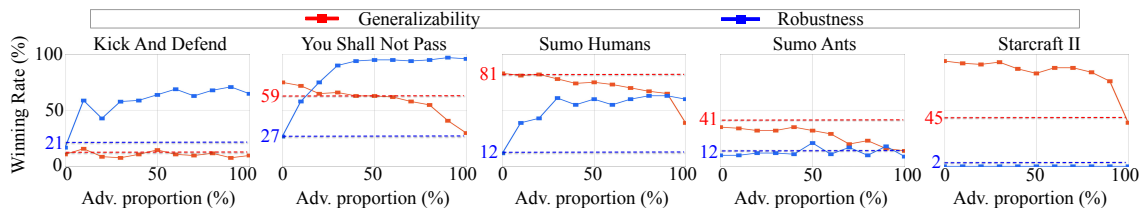
	Kick And Defend	You Shall Not Pass	Sumo Humans	Sumo Ants	StarCraft II
Mean/Std (%)	8.2/2.8	24.1/9.4	14.6/9.0	37.5/6.2 (Non-loss rate)	56.9/13.2
P-value	0.002	0.003	0.007	<0.001 (Non-loss rate)	0.002

a statistical measure on the winning rate difference between our attack ( $s_o$ ) and the baseline ( $s_b$ ). Specifically, we first compute their diff ( $D = s_o - s_b$ ) in each run. Then, we compute the mean, std, and the  $p$ -value of the paired t-test. For the t-test, our null hypothesis is  $H_0 : \mathbb{E}[D] \leq 0$ . If  $p$ -value is larger than an empirical threshold (*e.g.*, 0.01), we accept  $H_0$ , indicating our method cannot outperform the baseline. The results in Table S2 indicate a rejection of  $H_0$ , confirming our method’s superiority over the baseline approach.

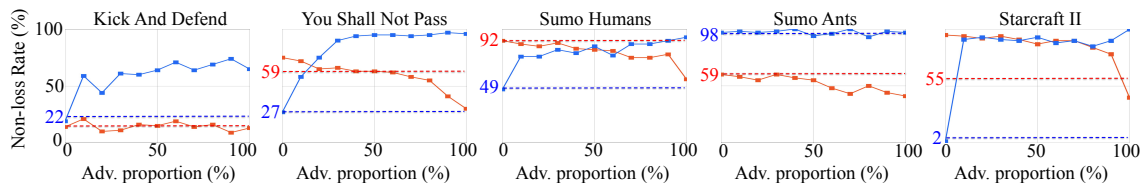
**Our attack vs. other comparison methods.** In this experiment, we compare our attack with two other methods of training an adversarial agent. We first consider an adversarial learning algorithm that optimizes a similar objective with our attack but without monotonic property. As is discussed in Section 3.2, a trivial way to solve the our proposed objective function is to approximate the second term in our learning objective with a DNN  $G_{\pi^\alpha}^v(\cdot)$  and combine it with the TRPO objective function (*i.e.*,  $M_{\pi^\alpha}^\alpha(\cdot)$ ). By maximizing this approximation of the objective function *i.e.*,  $M_{\pi^\alpha}^\alpha(\cdot) - G_{\pi^\alpha}^v(\cdot)$  with stochastic gradient ascent, one could solve an adversarial policy that shares a similar learning goal with our attack. However, this method cannot guarantee a monotonical increase in the difference between the expected total rewards of the adversarial agent and the victim agent. In this experiment, we apply this trivial solution to the selected games and compare the performance of the trained agent with the adversarial agent prepared by our attack.

We also compare our attack with another method that utilizes only the second term in our final learning objective function as the objective function (*i.e.*,  $\text{argmax}_\theta - \mathbb{E}[\min(\text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon)A_t^v, \rho_t A_t^v)]$ ). Recall that our final learning objective function contains two objectives. The first is to maximize the expected total reward of the adversarial agent, whereas the second is to minimize that of the victim. By setting up this experiment, we study the effect of the second objective upon our attack against the victim because intuition suggests the minimization of victim reward alone could also downgrade the performance of the victim and thus lead the potential victory of the adversary.

While this objective is distinct from that of (Gleave et.al 2020), it is my impression that the approach of (Gleave et.al 2020) is subtly mischaracterized in that it was meant to train policies which are adversarial to a particular policy (in that they minimize the victim’s reward), and not self interested (in the sense that they maximize their own reward). In the zero-sum environments these correspond very closely, but in the Starcraft II environment it seems that the suitable baseline would not be to train on the adversary’s environment reward, but on the victim’s environment reward. This seems to be a reasonable model of “adversarial behavior” even in general-sum games, but many applications it would be unrealistic (for instance adversarial cars following this model would crash into the victim with no regards to their own safety). Changing this aspect of the paper not only would be a more accurate reflection of prior work, but would likely clarify the difference between this approach and prior work. To further clarify this difference I would suggest to give the new objective a name that is not “adversarial” as it does not follow the typical usage of that framing. You use “hostile” at one point, which could be a suitable replacement as it does not imply that the agent is directly in opposition to the victim.



(a) The winning rates of the retrained victim agents.



(b) The non-loss rates of the retrained victim agents.

Figure S4: The performance of the victim agents retrained with different proportions of normal/adversarial episodes in the retraining episodes. The solid blue lines represent the average winning (plus tie) rates when playing with our adversarial agent (indicated by “robustness”). The solid red lines represent the average winning (plus tie) rates when playing with the second-strongest regular agent (indicated by “generalizability”). The blue and red dash lines represent the robustness and generalizability of the victim agent before the adversarial retraining, respectively. The x-axis label “Adv. proportion” denotes the proportions of adversarial episodes.

We conducted an experiment making adversary focus on minimizing victim reward without caring its own. Supplementary S6 shows the results (minimization attack in Fig. S2). We found, sometimes, this minimization method works but can’t outperform ours. For StarCraft II, the minimization method completely fails. This aligns with the reviewer’s thoughts. We will emphasize this point in the next version.

Figure S2 shows the performance of the adversarial agent obtained by our attack and the two comparison methods introduced above. In Figure S2, we can first observe that, without the monotonic guarantee, the adversarial winning rate dramatically drops on all of the five games. In all the games except “You-Should-Not-Pass”, the attack without monotonicity performs even worse than the state-of-art attack [4]. This result shows that, without the monotonic guarantee, the newly added minimization term often imposes even a negative impact upon the adversarial learning process and thus result in an adversarial policy with a weaker exploitability than that obtained by the attack without the minimization term. It should be noted that, in the complicated StarCraft II game, the learning process completely fails, resulting in zero winning plus tie rate. This indicates that enabling a monotonic guarantee is especially crucial for sophisticated games.

As we can also observe from Figure S2, by taking only the minimization into account alone, the adversarial agents trained do not demonstrate a similar winning rate as our proposed method. However, as is shown in Figure 2(b), they generally exhibit a powerful ability to prevent the victim from winning the game. As we can observe from the games “Sumo Ants and Humans”, the ability to downgrade the winning rate of the victim is almost as same as or even better than our proposed method. It indicates that the expected reward minimization could play a critical role in restricting the victim agent’s win.

However, from Figure S2, we can also observe the minimization alone does not make any significant difference for StarCraft game. It can barely bring any game wins or ties. We believe the reason is the design of the game. At the beginning of the StarCraft game, the game engine starts both agents at two different corners on a large map. As such, an adversarial agent in the StarCraft game cannot quickly interact with the victim, influence its actions, and thus curtail the victim’s reward collection. Instead, to prevent the victim from collecting resources and building up a strong army, the adversarial agent has to first spend time exploring the map and navigating to the base station of the victim. In this period, the victim usually has already constructed an army which is sufficiently strong to beat

intruders. Using our approach, which combines both minimizing victim’s reward and maximizing the adversarial reward, we can train an agent to have the ability first to gather resources to build up an army and then intervene in the army construction of the victim. With this learned policy, we can eventually obtain decisive wins.

Another point regarding the minimization attack is that, compared to the baseline approach, the policy trained by this attack is more likely to be *adversarial*. This is because rather than being self-interested by maximizing its own reward, the adversarial agent focuses on disturbing the victim agent. However, this attack can be unrealistic in many applications where self-interest is essential for obtaining a feasible policy (*e.g.*, StarCraft and adversarial self-driving cars – following this model would crash into the victim without considering their own safety).

**Adversarial retraining with different episode splits.** In Section 4.2, we follow the setup in [4] and set the adversarial and normal episodes evenly when conducting the adversarial retraining against our attack. In this experiment, we study the influence of episode split upon the retraining performance. Specifically, we vary the proportion of the adversarial episodes from 0.0%  $\sim$  100.0% by increasing 10% each time and retrain the victim agent by using each episode split. Figure S4 shows the robustness and the generalizability of the victim agent retrained under these settings. As we can first observe from the figure, overall, a higher proportion of adversarial episode enables better robustness, but worse generalizability. As is also shown in the figure, adversarial retraining has different effects on different games. For example, in You-Should-Not-Pass and Kick-And-Defend, the adversarial retraining always improves the robustness of the retrained victim agent, even with only 10% of adversarial episodes. It should be noted that, if one intends to select a setting, where both the robustness and generalizability are improved, different games will give different results. This indicates there is no universal optimal episode split, and the user has to find the best solution for each game individually. It should also be noted that, in Sumo-Ants, no matter how to vary the episode split, the robustness of the retrained victim keeps almost unchanged, which means adversarial retraining cannot robustify the victim agent against our attack. This result supports that, in this game, our adversarial policy exploits the game unfairness rather than the weakness of the opponent policy.

**Adversarial Policy Behavior Analysis.**

In addition to drawing the demo videos, we also follow [4] and conducts two additional experiments on the selected games to analyze the behavior of our adversarial policies. Specifically, we first mask a victim agent (*i.e.*, zero out the part of the victim observation that corresponds to the adversarial position) and play it with our adversarial agent and that obtained by the existing attack in the corresponding game. Note that the adversarial agents are trained against the original unmasked victims. Table S3 records the victim’s winning and non-loss rate before and after masking. As we can first observe from the Table, similar to the findings in [4], the victim winning rates against both attacks increase dramatically in the You-Should-Not-Pass and Kick-And-Defend game. This result also matches our observations from demo videos that an adversarial agent fails a victim by triggering adversarial observations rather than winning in a regular way. However, as is shown in Section 4, despite establishing similar adversarial behaviors, our adversarial agent has a stronger exploitability than that obtained by the

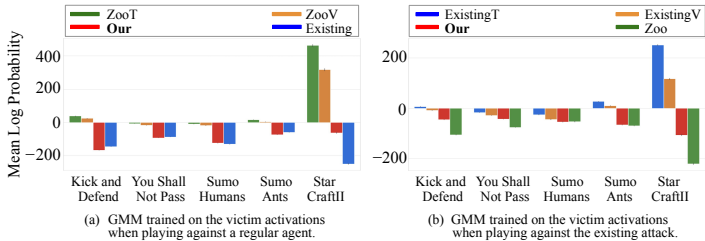


Figure S5: The average GMM log likelihood of the victim activations when playing against different opponents. Figure (a) shows the results of using the victim activations collected by playing against a regular agent to train the GMM. “ZooT” and “ZooV” represents the training/validation activations. “Our” and “Existing” represents the activations collected when playing against our attack and the existing attack. Figure (b) shows the results of training the GMM with the victim activations collected by playing against the existing attack (“ExistingT” and “ExistingV” are the training and testing set; “Our” and “Zoo” are our attack and regular agent).

Table S3: The victim winning/non-loss rates against our adversarial agents before and after masking.

			Kick And Defend	You Shall Not Pass	Sumo Humans	SumoAnts	StarCraft II
Our attack	Before masking	Winning (%)	14.0	26.0	22.0	15.0	1.0
		Non-loss (%)	14.0	26.0	53.0	86.0	2.0
	After masking	Winning (%)	94.0	98.0	25.0	13.0	4.0
		Non-loss (%)	94.0	98.0	67.0	87.0	24.0
Existing attack	Before masking	Winning (%)	45.0	48.0	34.0	57.0	18.0
		Non-loss (%)	45.0	48.0	65.0	91.0	64.0
	After masking	Winning (%)	97.0	97.0	12.0	37.0	65.0
		Non-loss (%)	97.0	97.0	68.0	95.0	98.0

Table S4: The winning rates of our adversarial agents against the mediocre and well-trained victims.

		Kick And Defend	You Shall Not Pass	Sumo Humans	Sumo Ants	StarCraft II
Adv. winning rate	Against mediocre victims (%)	88.0	93.0	55.0	86.0 (Non-loss rate)	98.0
	Against well-trained victims (%)	70.0	77.0	50.0	84.0 (Non-loss rate)	90.0

state-of-art approach. We also notice that masking almost has no impact upon our attack in the Sumo-Ants game, which further confirms that our adversarial agent exploits the game unfairness in this game. In addition, our attack establishes lower adversarial winning rate drop than the existing attack on the Starcraft II games. This confirms that our attack has a stronger exploitability than the existing attack on this sophisticated game. This may also indicate that our attack fails a victim via a stronger policy rather than triggering adversarial observations. Last but not least, we notice that masking even increases the adversarial winning rate of the existing attack on Sumo-Humans and Sumo-Ants. In [4], Gleave *et al.* also has the similar observation in their experiments. We suspect this is caused by the specific game rules of these two games. In future work, we will take a more closer look into the game rules together with the agent behaviors and find out the reasons behind this result.

Second, we collect the activations of the victim policy when playing with three different opponents: itself (a regular agent), our adversarial agent, and the adversarial agent obtained by the existing attack. We then follow the setup in [4] and use GMM and t-sne to demonstrate the differences among these sets of activations. Regarding GMM, we train two models with the activations collected from self-playing and playing against the existing attack. Then, we test these two models with these three sets of activations. The results are shown in Fig. S5 and Fig. S8. As we can first observe from these figures, on MuJoCo games, the results are aligned with our observations from the demo videos. That is, except for Sumo-Ants, our method exploits the similar weakness in victim policies with the existing attack. Regarding the StarCraft II game, Fig. S5 and Fig. S8 shows that the victim agent demonstrates substantially different behaviors when playing against our adversarial agent and that obtained by the existing attack. Together with the attack performance in Fig. 1 and Table S3, these results indicate our attack obtains an adversarial policy that exploits different and more threatening weakness than the existing attack on this sophisticated game.

**Attacking a mediocre victim.** As is shown in Supplementary S5.2, the victim agents used in our experiments are well-trained agents. Here, we show that our attack could also demonstrate its effectiveness even if we train our adversarial agents against mediocre agents. Specifically, We first train a mediocre agent on each game by running fewer self-play iterations. The average winning rate of these agents against the well-trained victims is 18%, confirming their mediocre performances. Then, we train our adversarial agent against these agents and test it against the mediocre victim and the well-trained victim. Table S4 shows the attack performances. We observe that our attack is effective against mediocre & well-trained victims even if the adversary is pitched on mediocre.

**Attacking a victim that varies its policy.** In our evaluation, we fix the victim agents. Here, we conduct an initial exploration of our attack’s effectiveness against a victim that varies its policy. Specifically, we train our attack with a victim that encodes two well-trained self-play policies and plays each one with an equal probability in each game round. Fig S6 shows the attack performance. The result confirms our attack’s effectiveness against this dynamic victim. Our future works will test more non-fixed victims.

Table S5: The adversarial agent’s non-loss rate against the victim agent and another regular agent.

		Kick And Defend	You Shall Not Pass	Sumo Humans	Sumo Ants	StarCraft II
Our attack	Victim (%)	91.0	89.0	94.0	89.0	98.0
	Regular agent (%)	70.0	60.0	82.0	88.0	92.0
Baseline attack	Victim (%)	83.0	84.0	93.0	61.0	64.0
	Regular agent (%)	60.0	64.0	74.0	45.0	58.0

## S7 Attack Transferability

We also compare the transferability of our attack with that of the baseline attack [4]. Specifically, given a game, we first take a regular agent released in [2]. This agent is different from the one used for adversarial policy training. Then, we set up the regular agent to play with the adversarial agent learned through our method and that learned through the baseline approach [4]. In this experiment, we set each adversarial agent to play with the regular agent for 100 rounds and report the adversarial agent’s winning rate. We compare the adversary’s winning rate with the winning rate observed when the adversary plays with the victim it trains against. Through this comparison, we can measure the exploitability variation after transferring an adversarial agent to attack a different target agent. Recall that for each game, we randomly select six initial states and thus obtain six adversarial agents. In this experiment, we report the transferability of the strongest adversarial agent. Table S5 shows the transferability of our attack and that of the baseline attack [4]. First, we observe that both methods establish a certain level of transferability on the five games. Compared with the baseline attack, our attack demonstrates a slightly better transferability. We believe this is because of the stronger exploitability of our attack. As is also shown in the table, our adversarial policy in Sumo-Ants establishes the highest transferability. As is mentioned above, this policy exploits the game unfairness rather than the weakness of a specific victim policy. As such, it is less relevant to the opponent policy and thus has a stronger transferability. On the contrary, as for the adversarial policy that disturbs the victim observation via its action, its performance will be jeopardized when transferred to a different regular policy.

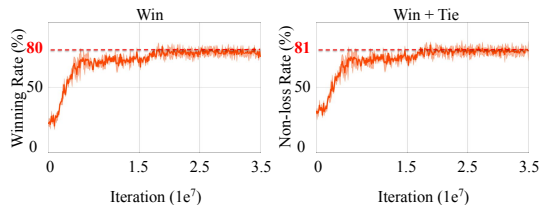


Figure S6: Our attack performance against a dynamic victim in the Kick And Defend game.

## S8 Our attack vs. baseline in zero-sum setting

In Section 4, we show the advantages of our attack over the baseline [4] on real-world nonzero-sum games. Here, we demonstrate the performance of both methods in a zero-sum setting. Specifically, we first modify the reward design of the StarCraft II and make it a zero-sum game.<sup>2</sup> To achieve this, we remove all the intermediate rewards and only preserve the rewards related to the game results, i.e., 1 (win), 0 (tie), and -1 (lose). We then apply both our attack and the baseline to train an adversarial agent under the modified reward design and record the adversarial winning rate every time its policy is updated. Figure S7 shows the results of six runs. As we can observe from the figure, the adversarial agent trained by the baseline [4] demonstrates a similar winning rate with our attack. In zero-sum games, maximizing adversarial reward will automatically minimize

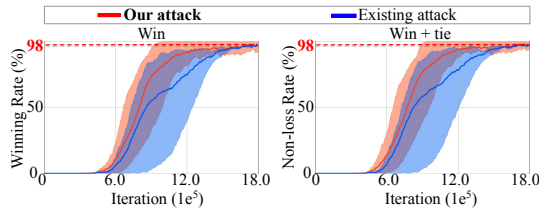


Figure S7: Performance comparison in a zero-sum game.

<sup>2</sup>The MuJoCo games cannot be transformed into zero-sum settings, because the agent is not able to pick up basic behaviors, such as standing and running, without the intermediate rewards.

the victim reward. As such, our approach’s adversarial agent will demonstrate the same performance as that learned by the baseline [4] in this setting. It should be noted that the adversarial agent trained by the baseline in the zero-sum setting performs better than that in the nonzero-sum setting (Fig. 1). As is discussed above, this is because the state-of-art attack is less effective in nonzero-sum settings. It should also be noted that our learning converges much faster in the nonzero-sum setting (Fig. 1: 0.8M iterations) than it in the zero-sum setting (1.6M), which demonstrates the benefit brought by the intermediate rewards. In most of the two-player games, game developers design intermediate rewards to help RL agent training.

## References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *Proc. of OSDI*, 2016.
- [2] Trapit Bansal, Jakub Pachocki, Szymon Sidor, Ilya Sutskever, and Igor Mordatch. Emergent complexity via multi-agent competition. In *Proc. of ICLR*, 2018.
- [3] DeepMind. Alphastar: Mastering the real-time strategy game starcraft ii. [https://en.wikipedia.org/wiki/AlphaStar\\_\(software\)](https://en.wikipedia.org/wiki/AlphaStar_(software)), 2017.
- [4] Adam Gleave, Michael Dennis, Neel Kant, Cody Wild, Sergey Levine, and Stuart Russell. Adversarial policies: Attacking deep reinforcement learning. In *Proc. of ICLR*, 2020.
- [5] David A Levin and Yuval Peres. *Markov chains and mixing times*. American Mathematical Society., 2017.
- [6] Ruo-Ze Liu, Haifeng Guo, Xiaozhong Ji, Yang Yu, Zhen-Jia Pang, Zitai Xiao, Yuzhou Wu, and Tong Lu. Efficient reinforcement learning with a thought-game for starcraft. *arXiv preprint arXiv:1903.00715*, 2019.
- [7] openai. Stable baselines. <https://stable-baselines.readthedocs.io/en/master/>, 2018.
- [8] Jan Peters, Katharina Mulling, and Yasemin Altun. Relative entropy policy search. In *Proc. of AAAI*, 2010.
- [9] Jan Peters and Stefan Schaal. Natural actor-critic. *Neurocomputing*, 2008.
- [10] David Pollard. Asymptopia: an exposition of statistical asymptotic theory. 2000. URL <http://www.stat.yale.edu/pollard/Books/Asymptopia>, 2000.
- [11] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *Proc. of ICML*, 2015.
- [12] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [13] Peng Sun, Xinghai Sun, Lei Han, Jiechao Xiong, Qing Wang, Bo Li, Yang Zheng, Ji Liu, Yongsheng Liu, Han Liu, et al. Tstarbots: Defeating the cheating level builtin ai in starcraft ii in the full game. *arXiv preprint arXiv:1809.07193*, 2018.
- [14] Sebastian Thrun. Monte carlo pomdps. In *Proc. of NeurIPS*, 2000.
- [15] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *Proc. of ICIRS*, 2012.
- [16] Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.

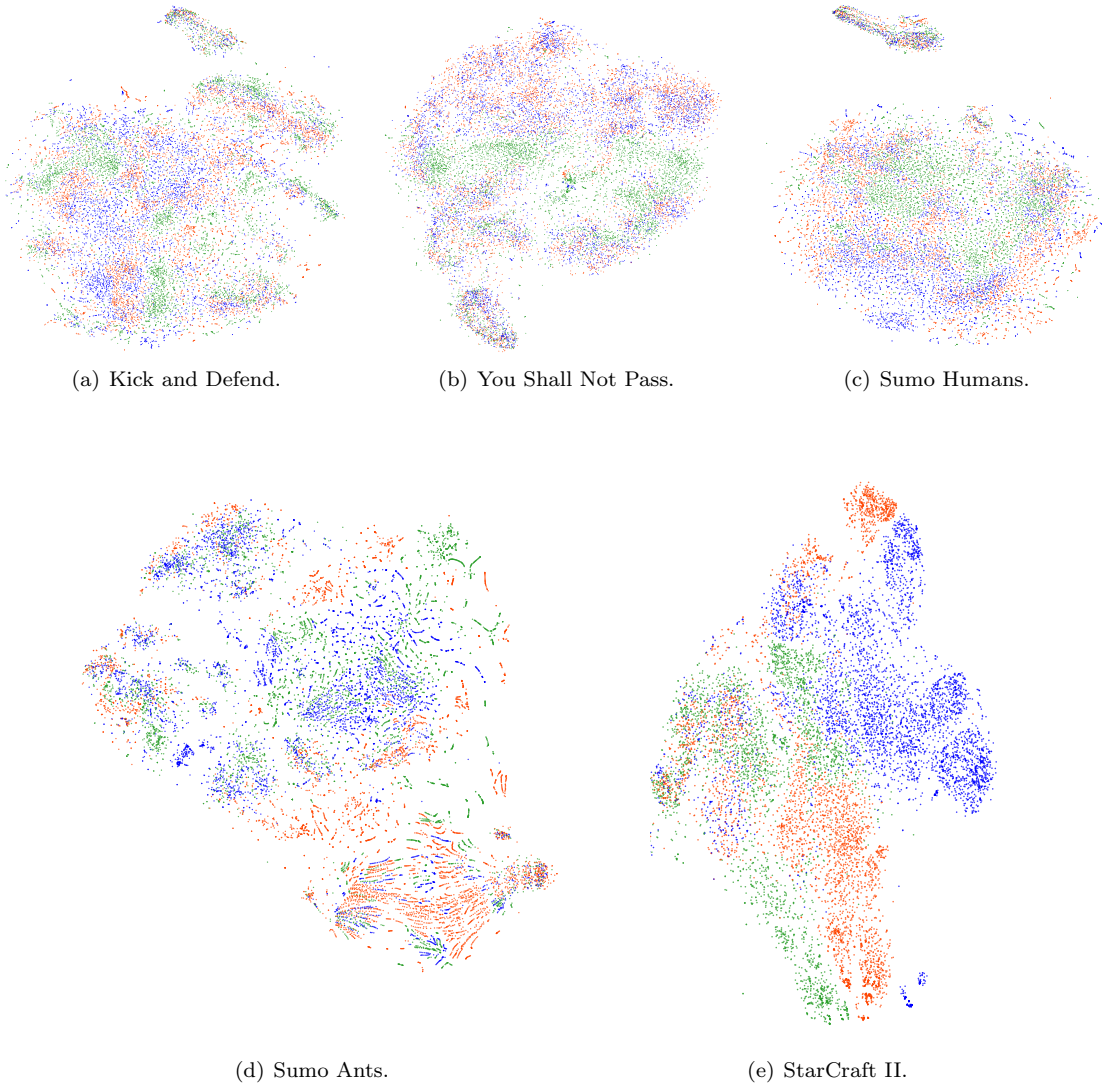


Figure S8: t-SNE visualizations of the victim activations when playing against different opponents in MuJoCo games. The green dots are the victim activations when setting a regular agent as the opponent. The red and blue dots indicates the victim activations when playing against our adversarial agent and that obtained by the existing attack, respectively.