# Towards Cooperation in Sequential Prisoner's Dilemmas: a Deep Multiagent Reinforcement Learning Approach

**Weixun Wang**[1]**, Jianye Hao** [1]**, Yixi Wang** [1]**, Matthew Taylor**[2]**,**

[1] Tianjin University, Tianjin, China

[2] Washington State University, Pullman, WA, USA

wxwang@tju.edu.cn, jianye.hao@tju.edu.cn, yixiwang2017@outlook.com, taylorm@eecs.wsu.edu

## Abstract

The Iterated Prisoner's Dilemma has guided research on social dilemmas for decades. However, it distinguishes between only two atomic actions: cooperate and defect. In real-world prisoner's dilemmas, these choices are temporally extended and different strategies may correspond to sequences of actions, reflecting grades of cooperation. We introduce a Sequential Prisoner's Dilemma (SPD) game to better capture the aforementioned characteristics. In this work, we propose a deep multiagent reinforcement learning approach that investigates the evolution of mutual cooperation in SPD games. Our approach consists of two phases. The first phase is offline: it synthesizes policies with different cooperation degrees and then trains a cooperation degree detection network. The second phase is online: an agent adaptively selects its policy based on the detected degree of opponent cooperation. The effectiveness of our approach is demonstrated in two representative SPD 2D games: the Apple-Pear game and the Fruit Gathering game. Experimental results show that our strategy can avoid being exploited by exploitative opponents and achieve cooperation with cooperative opponents.

## 1 Introduction

Learning is the key to achieving coordination with others in multiagent environments [Stone and Veloso, 2000]. Over the last couple of decades, a large body of multiagent learning techniques have been proposed that aim to coordinate on various solutions (e.g., Nash equilibrium) in different settings, e.g., minimax Q-learning [Littman, 1994], Nash Q-learning [Hu and Wellman, 2003], and Conditional-JAL [Banerjee and Sen, 2007], to name just a few.

One commonly investigated class of games is the Prisoner's Dilemma (PD), in which an Nash equilibrium solution is not a desirable learning target. Until now, a large body of work [Axelrod, 1984; Nowak and Sigmund, 1993; Banerjee and Sen, 2007; Crandall and Goodrich, 2005; Damer and Gini, 2008; Hao and Leung, 2015; Mathieu and Delahaye, 2015] has been devoted to incentivize rational agents towards mutual cooperation in repeated matrix PD games. However, all the above works focus on the classic repeated PD games, which ignores several key aspects of real world prisoner's dilemma scenarios. In repeated PD games, the moves are atomic actions and can be easily labeled as cooperative or uncooperative or learned from the payoffs [Busoniu *et al.*, 2008]. In contrast, in real world PD scenarios, cooperation/defection behaviors are temporally extended and the payoff signals are usually delayed (available after a number of steps of interactions).

Crandall [2012] proposes the Pepper framework for repeated stochastic PD games (e.g., two-player gate entering problem), which can extend strategies originally proposed for classic repeated matrix games. Later some techniques extend Pepper in different scenarios, e.g., stochastic games with a large state space under tabular based framework [Elidrisi *et al.*, 2014] and playing against the switching opponents [Hernandez-Leal and Kaisers, 2017]. However, these approaches rely on hand-crafted state inputs and tabular Q-learning techniques to learn optimal policies. Thus, they cannot be directly applied to more realistic environments whose states are too large and complex to be analyzed beforehand.

Leibo et al. [2017] introduce a 2D Fruit Gathering game to better capture the real world social dilemma characteristics, while also maintaining the characteristics of classical iterated PD games. In this game, at each time step, an agent selects its action based on its image observation and cannot directly observe the actions of the opponent. Different policies represent different levels of cooperativeness, which is a graded quantity. They investigate the cooperation/defection emergence problem by leveraging the power of deep reinforcement learning [Mnih *et al.*, 2013; Mnih *et al.*, 2015] from the descriptive point of view: how do multiple selfish independent agents' behaviors evolve when agents update their policy using deep Q-learning? In contrast, this paper takes a prescriptive and non-cooperative perspective and considers the following question: *how should an agent learn effectively in real world social dilemma environments when it is faced with different opponents?*

To this end, in the paper, we first formally introduce the general notion of sequential prisoner's dilemma (SPD) to model real world PD problems. We propose a multiagent deep reinforcement learning approach for mutual cooperation in SPD games. Our approach consists of two phases: offline and online phases. The offline phase generates policies with

varying cooperation degrees and trains a cooperation degree detection network. To generate policies, we propose using the weighted target reward and two schemes, IAC and JAC, to train the baseline policies with varying cooperation degrees. Then we propose using the policy generation approach to synthesize the full range of policies from these baseline policies. Lastly, we propose a cooperation degree detection network implemented as an LSTM-based structure with an encoder-decoder module and generate a training dataset. The online phase extends the Tit-for-Tat principle into sequential prisoner's dilemma scenarios. Our strategy adaptively selects its policy with the proper cooperation degree from a continuous range of candidates based on the detected cooperation degree of an opponent. Intuitively, on one hand, our overall algorithm is cooperation-oriented and seeks for mutual cooperation whenever possible; on the other hand, our algorithm is also robust against selfish exploitation and resorts to a defection strategy to avoid being exploited whenever necessary. We evaluate the performance of our deep multiagent reinforcement approach using two 2D SPD games (the Fruit Gathering and Apple-Pear games). Our experiments show that our agent can efficiently achieve mutual cooperation under self-play and also perform well against opponents with changing stationary policies.

## 2 Background

### 2.1 Matrix Games and the Prisoner's Dilemma

A matrix game can be represented as a tuple $\langle N, \{A_i\}_{i \in N}, \{R_i\}_{i \in N}\rangle$, where $N$ is the set of agents, $A_i$ is the set of actions available to agent $i$ with $\mathcal{A}$ being the joint action space $A_1 \times \ldots \times A_n$, and $R_i$ is the reward function for agent $i$. One representative class of matrix games is the prisoner's dilemma game, as shown in Table 1. In this game, each agent has two actions: cooperate (C) and defect (D), and is faced with four possible rewards: $R$, $P$, $S$, and $T$. The four payoffs satisfy the following four inequalities under a prisoner's dilemma game:

- $R > P$: mutual cooperation is preferred to mutual defection.

- $R > S$: mutual cooperation is preferred to being exploited by a defector.

- $2R > S+T$: mutual cooperation is preferred to an equal probability of unilateral cooperation and defection.

- $T > R$: exploiting a cooperator is preferred over mutual cooperation.

- $P > S$: mutual defection is preferred over being exploited.

### 2.2 Markov Game

Markov games combine matrix games and Markov Decision Processes and can be considered as an extension of Matrix games to multiple states. A Markov game $\mathcal{M}$ is defined by a tuple $\langle N, S, \{A_i\}_{i \in N}, \{R_i\}_{i \in N}, \mathcal{T}\rangle$, where $S$ is the set of states and $N$ is the number of agents, $\{A_i\}_{i \in N}$ is the collection of action sets, with $A_i$ being the action set of agent $i$, and $\{R_i\}_{i \in N}$ is the set of reward functions,

Table 1: Prisoner's Dilemma

|   | C | D |
|---|---|---|
| C | R, R | S, T |
| D | T, S | P, P |

$R_i : S \times A_1 \times \ldots \times A_n \to \mathcal{R}$ is the reward function for agent $i$. $\mathcal{T}$ is the state transition function: $S \times A_1 \times \ldots \times A_n \to \Delta(S)$, where $\Delta(S)$ denotes the set of discrete probability distributions over $S$. Matrix games are the special case of Markov games when $|S| = 1$.

Next we formally introduce SPD by extending the classic iterated PD game to multiple states.

### 2.3 Definition of Sequential Prisoner's Dilemma

A two-player SPD is a tuple $\langle \mathcal{M}, \Pi\rangle$, where $\mathcal{M}$ is a 2-player Markov game with state space $S$. $\Pi$ is the set of policies with varying cooperation degrees. The empirical payoff matrix $(R(s), P(s), S(s), T(s))$ can be induced by policies $(\pi^C, \pi^D \in \Pi)$, where $\pi^C$ is more cooperative than $\pi^D$. Given two policies, $\pi^C$ and $\pi^D$, the corresponding empirical payoffs $(R, P, S, T)$ under any starting state s with respect to the payoff matrix in Section 2.1 can be defined as $(R(s), P(s), S(s), T(s))$ through their long-term expected payoff, where

$$R(s) := V_1^{\pi^C, \pi^C}(s) = V_2^{\pi^C, \pi^C}(s), \qquad (1)$$

$$P(s) := V_1^{\pi^D, \pi^D}(s) = V_2^{\pi^D, \pi^D}(s), \qquad (2)$$

$$S(s) := V_1^{\pi^C, \pi^D}(s) = V_2^{\pi^D, \pi^C}(s), \qquad (3)$$

$$T(s) := V_1^{\pi^D, \pi^C}(s) = V_2^{\pi^C, \pi^D}(s), \qquad (4)$$

We can define the long-term payoff $V_i^{\vec{\pi}}(s)$ for agent $i$ when the joint policy $\vec{\pi} = (\pi_1, \pi_2)$ is followed starting from $s$.

$$V_i^{\vec{\pi}}(s) = E_{\vec{a}_t \sim \vec{\pi}(s_t), s_{t+1} \sim \mathcal{T}(s_t, \vec{a}_t)}[\sum_{t=0}^{\infty} \gamma^t r_i(s_t, \vec{a}_t)] \qquad (5)$$

A Markov game is an SPD when there exists a state $s \in S$ for which the induced empirical payoff matrix satisfies the five inequalities in Section 2.1. Since SPD is more complex than PD, the existing approaches addressing learning in matrix PD games cannot be directly applied in SPD.

### 2.4 Deep Reinforcement Learning

**Q-Learning and Deep Q-Networks**: Q-learning and Deep Q-Networks (DQN) [Mnih *et al.*, 2013; Mnih *et al.*, 2015] are value-based reinforcement learning approaches to learn optimal policies in Markov environments. Q-learning makes use of an action-value function for policy $\pi$ as $Q^{\pi}(s, a) = E_{s'}[r(s, a) + \gamma E_{a' \sim \pi}[Q^{\pi}(s', a')]]$. DQN uses a deep convolutional neural network to estimate Q-values and the optimal Q-values are learned by minimizing the following loss function:

$$y = r + \gamma max_{a'}\bar{Q}(s', a'), \qquad (6)$$

$$L(\theta) = E_{s,a,r,s'}[(Q(s, a|\theta) - y)^2] \qquad (7)$$

where $\bar{Q}$ is a target $Q$ network whose parameters are periodically updated with the most recent $\theta$.

**Policy Gradient and Actor-Critic Algorithms**: Policy Gradient methods are for a variety of RL tasks [Williams, 1992; Sutton *et al.*, 2000]. Their objective is to maximize $J(\theta) = E_{s \sim p^\pi, a \sim \pi_\theta}[R]$ by taking steps in the direction of $\bigtriangledown_\theta J(\theta)$, where

$$\bigtriangledown_\theta J(\theta) = E_{s \sim p^\pi, a \sim \pi_\theta}[\bigtriangledown_\theta log \pi_\theta(a|s) Q^\pi(s,a)] \quad (8)$$

where $p^\pi$ is the state transition distribution. Practically, the value of $Q^\pi$ can be estimated in different ways. For example, $Q^\pi(s,a)$ serves as a critic to guide the updating direction of $\pi_\theta$, which leads to a class of actor-critic algorithms [Schulman *et al.*, 2015; Wang *et al.*, 2016].

## 3  Deep RL: Towards Mutual Cooperation

Algorithm 1 describes our deep multiagent reinforcement learning approach, which consists of two phases, as discussed Section 1. In the offline phase, we first seek to generate policies with varying cooperation degrees. Since the number of policies with different cooperation degrees is infinite, it is computationally infeasible to train all the policies from scratch. To address this issue, we first train representative policies using Actor-Critic until convergence (i.e., cooperation and defection baseline policies) (Lines 3-5) detailed in Section 3.1; second, we synthesize the full range of policies (Lines 6-7) from the above baseline policies, which will be detailed in Section 3.2. Another task is to how to effectively detect the cooperation degree of the opponent. We divide this task into two steps: we first train an LSTM-based cooperation degree detection network offline (Lines 8-10), which will be then used for real-time detection during the online phase, detailed in Section 3.3. In the online phase, our agent plays against any opponent by reciprocating with a policy of a slightly higher cooperation degree than that of the opponent we detect (Lines 12-18), detailed in Section 3.4. Intuitively, on one hand, our algorithm is cooperation-oriented and seeks for mutual cooperation whenever possible; on the other hand, our algorithm is also robust against selfish exploitation and resorts to defection strategy to avoid being exploited whenever necessary.

### 3.1  Train Baseline Policies with Different Cooperation Degrees

One way of generating policies with different cooperation degrees is by directly changing the key parameters of the environments. For example, Leibo et al. [Leibo *et al.*, 2017] investigate the influence of the resource abundance degree on the learned policy's cooperation tendency in sequential social dilemma games where agents compete for limited resources. It is found that when both agents employ deep Q-learning algorithms, more cooperative behaviors can be learned when resources are plentiful and vice versa. We may leverage similar ideas of modifying game settings to generate policies with different cooperation degrees. However, this type of approach requires a perfect understanding of the environment as a prior, and also may not be practically feasible when on cannot modify the underlying game engine.

Another more generalized way of generating policies with different cooperation degrees is to modify agents' reward signals during learning. Intuitively agents with the reward of the

---

**Algorithm 1** The Approach of Deep Multiagent Reinforcement Learning Towards Mutual Cooperation

---

1: *//offline training*
2: initialize the size $N_t$ of training policy set, the size $N_g$ of generation policy set, the number $N_d$ of training data set, the episode number $N_e$ and the step number $N_r$ of each episode
3: **for** training policy set index t = 1 to $N_t$ **do**
4:     set agents' attitudes
5:     train agents' policy set $P_t$ using weighted target reward
6: **for** generation policy set index g = 1 to $N_g$ **do**
7:     use policy set $\{P_t\}_{t \in N_t}$ to generate policy set $P_g$
8: **for** training data set index d = 1 to $N_d$ **do**
9:     generate training data set $D_d$ as $\{P_t\}_{t \in N_t} \cup \{P_g\}_{g \in N_g}$
10: use data set $\{D_d\}_{d \in N_d}$ to train cooperation degree detection network
11: *//adjust the policy online*
12: initialize $agent_1's$ cooperation degree
13: **for** episode index e = 1 to $N_e$ **do**
14:     **for** step index r = 1 to $N_r$ **do**
15:         $agent_1$ and $agent_2$ take actions and get rewards
16:         $agent_1$ uses $n$-state trajectory $\langle s_{r-n+1}, \ldots, s_r \rangle$ to detect the cooperation degree $cd_2^r$ of $agent_2$
17:         $agent_1$ updates $cd_1^r$ incrementally based on $cd_2^r$
18:         $agent_1$ synthesizes a policy with $cd_1^r$ using policy generation

---

sum of all agents' immediate rewards would learn towards cooperation policies to maximize the expected accumulated social welfare eventually, and agents maximizing only their own reward would learn more selfish (defecting) policies. Formally, for a two-player environment (agent $i$ and $j$), agent $i$ computes a weighted target reward $r_i'$ as follows:

$$r_i' = r_i + att_{ij} \times r_j \quad (9)$$

where $att_{ij} \in [0,1]$ is agent $i$'s attitude towards agent $j$, which reflects the relative importance of agent $j$ in agent $i's$ perceived reward. By setting the values of $att_{ij}$ and $att_{ji}$ to 0, agents would update their strategies in the direction of maximizing their own accumulated discounted rewards. By setting the values of $att_{ij}$ and $att_{ji}$ to 1, agents would update their strategies in the direction of maximizing the overall accumulated discounted reward. The greater the value of agent one's attitude towards agent two is, the higher the cooperation degree of agent one's learned policy would be.

Given the modified reward signal for each agent, the next question is how agents should learn to effectively converge to the expected behaviors. One natural way is to resort to the the independent Actor-Critic (IAC) or some other independent deep reinforcement learning, i.e., equipping each agent with an individual deep Q-learning algorithm (IDQL) with the modified reward. However, the key element to the success of DQL, experience replay memory, might prohibit effective learning in deep multiagent Q-learning environments [Foerster *et al.*, 2017b; Sunehag *et al.*, 2017]. The nonstationarity introduced by the coexistence of multiple IACs means that data in the replay memory may no longer reflect the current dynamics in which the agent is learning. Thus IACs may fre-
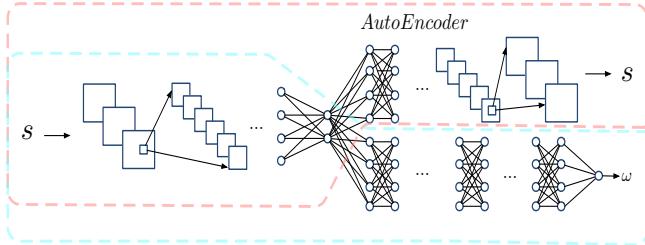
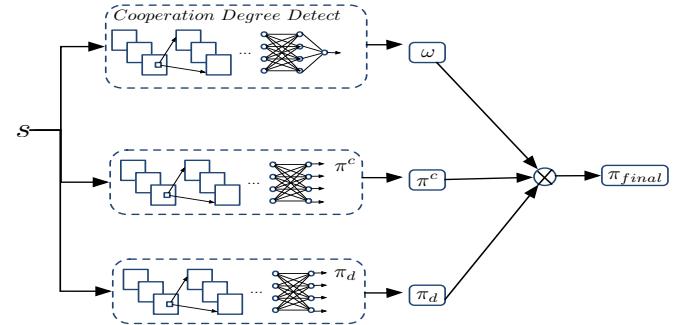Figure 1: Cooperation Degree Detection Network



Figure 2: The Structure of Deep Reinforcement Learning Approach Towards Mutual Cooperation



Figure 3: The Fruit Gathering game (left) and the Apple-Pear game.

quently get confused by obsolete experience and impede the learning process. A number of methods have been proposed to remedy this issue [Foerster *et al.*, 2017b; Lowe *et al.*, 2017; Foerster *et al.*, 2017a] and we omit the details which are out of the scope of this paper.

Since the baseline policy training step is performed offline, another way of improving training is to use the joint Actor-Critic (JAC) by treating both agents as a single learner during training. Note that we use JAC only for the purpose of training baseline policies offline, and we do not require that we can control the policies that an opponent may use online. In JAC, both agents share the same underlying network that learns the optimal policy over the joint action space using a single reward signal. In this way, the aforementioned nonstationarity problem can be avoided. Besides, compared with IAC, the training efficiency can be improved significantly since the network parameters are shared across agents in JAC.

In JAC, the weighted target reward is defined as follows:

$$r_{total} = \sum_{i \in N} att_i \times r_i \qquad (10)$$

where $att_i \in [0,1]$ represents the relative importance of agent $i$ on the overall reward. The smaller the value of $att_i$, the higher the cooperation degree of agent $i$'s learned policy and vice versa. Given the learned joint policy $\pi^{joint}(a_1, a_2|s, att_1, att_2)$, agent $i$ can easily obtain its individual policy $\pi_i$ as follows:

$$\pi_i = \sum_{a_j, \, j \neq i} \pi^{joint}(a_1, a_2|s, att_1, att_2) \qquad (11)$$

As we mentioned previously, it is computationally prohibitive to train a large number of policies with different cooperation degrees due to the high training cost of deep Q-learning, and because the policy space is infinite. To alleviate this issue, here we propose that only two policies, cooperation policy $\pi_c$ and defection policy $\pi_d$, need to be trained. Other policies with cooperation degree between the baselines can be synthesized efficiently, which will be introduced in Section 3.2.

### 3.2 Policy Generation

Given the baseline policies $\pi_i^c$ and $\pi_i^d$, we synthesize multiple policies. Each continuous weighting factor $w_c \in [0,1]$ corresponds to a new policy $\pi_i^{w_c}$ defined as follows:

$$\pi_i^{w_c} = w_c \times \pi_i^c + (1 - w_c) \times \pi_i^d \qquad (12)$$

The weighting factor $w_c$ is defined as policy $\pi_i^{w_c}$'s cooperation degree. Specially, the linear combination of two policies has two advantages — it 1) generates policies with varying cooperation degrees and 2) ensures low computational cost. Any synthesized policy $\pi_i^{w_c}$ should be more cooperative than $\pi_i^d$ and more defecting than $\pi_i^c$. The higher the value of $w_c$ is, the more cooperative the corresponding policy $\pi_i^{w_c}$ is. It is important to mention that the cooperation degrees of synthesized policies are ordinal, i.e., the cooperation degree of policies only reflect their relative cooperation ranking. For example, considering two synthesized policies $\pi_i^{0.6}$ and $\pi_i^{0.3}$, it only implies that policy $\pi_i^{0.6}$ is more cooperative than policy $\pi_i^{0.3}$. However, we cannot say that $\pi_i^{0.6}$ is twice as cooperative as $\pi_i^{0.3}$. Our way of synthesizing new policies can be understood as synthesizing policies over expert policies [He *et al.*, 2016]. The previous work applies a similar idea to generate policies to better respond with different opponents in competitive environments, however, our goal here is to synthesize policies with varying cooperation degrees in sequential Prisoner's dilemmas.

### 3.3 Opponent Cooperation Degree Detection

In classical iterated PD games, a number of techniques have been proposed to estimate the opponent's cooperation degree, e.g., counting the cooperation frequency when action can be observed; using a partial filter or a Bayesian approach otherwise [Damer and Gini, 2008; Hernandez-Leal *et al.*, 2016a; Hernandez-Leal *et al.*, 2016b; Leibo *et al.*, 2017]. However, in SPD, the actions are temporally extended, and the opponent's information (actions and rewads) cannot be observed
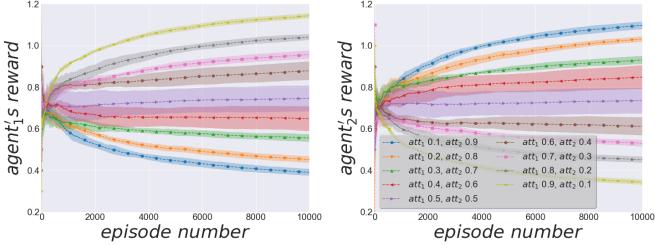
Figure 4: Agents' average rewards under policies trained with different cooperation attitudes.

directly. Thus the previous works cannot be directly applied. We need a way of accurately predicting the cooperation degree of the opponents from the observed sequence of moves in a qualitative manner. In SPD, given the sequential observations (time-series data), we propose an LSTM-based cooperation degree detection network.

In previous sections, we have introduced a way of synthesizing policies of any cooperation degree, thus we can easily prepare a large dataset of agents' behaviors with varying cooperation degrees. Based on this, we can transform the cooperation degree detection problem into a supervised learning problem: given a sequence of moves of an opponent, our task is to detect the cooperation degree (label) of this opponent. We propose a recurrent neural network, which combines an autoencoder and a recurrent classifier, as shown in Figure 1. Combing an autoencoder with a recurrent classifier brings two major benefits here. First, the classifier and autoencoder share underlying layer parameters of the neural network. This ensures that the classification task is based on the effective feature extraction of the observed moves, which improves the classification detection accuracy. Second, concurrent training of the autoencoder also helps to accelerate the training speed of the classifier and reduces fluctuation during training.

The network is trained on experiences collected by agents. Both agents $i$ and $j$ interact with the environment starting with initialized policies $\pi_i$ and $\pi_j$, yielding a training set $D$:

$$D = \{(X, label) | X = env(\pi_i, \pi_j),$$
$$\pi_i \in \{\pi_i^c, \pi_i^d\}, \pi_j \in \Pi_j, \qquad (13)$$
$$\text{if } \pi_i = \pi_i^c, \ label = 1 \text{ else } 0\}_{i \in N}$$

where $\pi_i^c$ and $\pi_i^d$ are baseline policies for agent $i$. $\Pi_j$ is the learned policy set of its opponent $j$. $label$ is the relative cooperation degree of policy $\pi_i$, and $X$ is the set of trajectories under the joint policy $(\pi_i, \pi_j)$. For each trajectory $x = \langle s_1, s_2, \ldots, s_d \rangle \in X$, its label is the cooperation degree of agent $i$'s policy. The network is trained to minimize the following weighted cross entropy loss function as follows:

$$L(f_c(\vec{x}), label) = -(w_1 \times label_1 \times log(p_1) +$$
$$w_2 \times (1 - lable_1) \times log(1 - p_1)) \qquad (14)$$

where $w_1$ is the weight of $label_1$. $p_1$ is the network output, which is the probability of $\pi^c$.

### 3.4 Play Against Different Opponents

Once we have detected the cooperation degree of the opponent, the final question arises as to how an agent should select proper policies to play against that opponent. A self-interested approach would be to simply play the best response policy toward the detected policy of the opponent. However, as we mentioned before, we seek a solution that can allow agents to achieve cooperation while avoiding being exploited.

Figure 2 shows our overall approach playing with opponents towards mutual cooperation. At each time step $t$, agent $i$ uses its previous $n$-step sequence of observations (from time step $t - n$ to $t$) as the input of the detection network, and obtain the detected cooperation degree $cd_j^t$. However, the one-shot detection of the opponent's cooperation degree might be misleading due to either the detection error of our classifier or the stochastic behaviors of the opponent. Thus this may lead to high variance of our detection outcome and our response policy thereafter. To reduce the variance, agent $i$ uses the exponential smoothing to update its current cooperation degree estimation of its opponent $j$ as follows:

$$cd_i^t = (1 - \alpha) \times cd_i^{t-1} + \alpha \times cd_j^t \qquad (15)$$

where $cd_i^{t-1}$ is agent $i$'s cooperation degree in the last time step, and $\alpha$ is the cooperation degree changeing factor.

Finally, agent $i$ sets its own cooperation degree equal to $cd_i^t$ plus its reciprocation level, and then synthesizes a new policy with the updated cooperation degree following Equation (12) as its next-step strategy to play against its opponent. Note that the reciprocation level can be quite low and still produce cooperation. The benefit is that it will not lead to a significant loss if the opponent is not cooperative at all, while full cooperation can be reached if the opponent reciprocates in the same way. Also, note that here we only provide a way of responding to the opponent with changing policies. However, our overall approach is general and any existing multiagent strategy selection approaches can be applied here.

## 4 Simulation and Results

### 4.1 SPD Game Descriptions

In this section, we adopt the Fruit Gathering game [Leibo *et al.*, 2017] to evaluate the effectiveness of our approach. We also propose another game Apple-Pear, which also satisfies real world social dilemma conditions we mentioned before. Each game involves two agents (in blue and red). The task of an agent in the Fruit Gathering game is to collect as many apples, represented by green pixels, as possible (see Figure 3 (left)). An agent's action set is: step forward, step backward, step left, step right, rotate left, rotate right, use beam, and stand still. The agent obtains the corresponding fruit when it steps on the same square as the fruit is located. When an agent collects an apple, it will receive a reward 1. And the apple will be removed from the environment and respawn after 40 frames. Each agent can also emit a beam in a straight line along its current orientation. An agent is removed from the map for 20 frames if it is hit by the beam twice. Intuitively, a defecting policy in this game is one which frequently tags the rival agent to remove it from the game. A cooperation

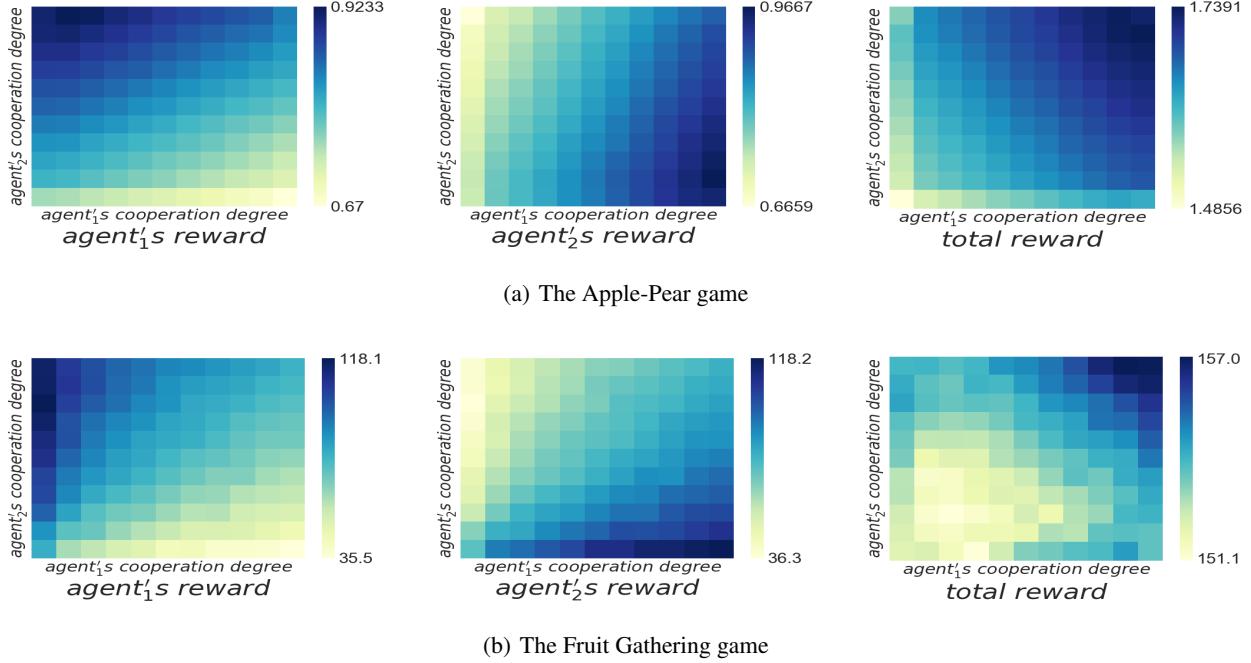(a) The Apple-Pear game



(b) The Fruit Gathering game

Figure 5: Average and total rewards under different cooperation degrees. The cooperation degrees of $agent_1$ and $agent_2$ increase from left to right and from bottom to top respectively. Each cell corresponds the rewards of different policy pairs.

policy is one that rarely tags the other agent. For the Apple-Pear game, there is a red apple and a green pear (see Figure 3 (right)). The blue agent prefers apple while the red agent prefers pear. Each agent has four actions: step right, step left, step backward, step forward, and each step of moving incurs a cost of 0.01. The fruit is collected when the agent steps on the same square as it. When the blue (red) agent collects an apple (pear) individually, it receives a higher reward 1. When the blue agent collects a pear individually, it receives a lower reward 0.5. The situation is the opposite for the red agent. One exception is that they both receive a half of their corresponding rewards when they share a pear or an apple. In this game, a fully defecting policy is to collect both fruits whenever the fruit-collecting reward exceeds the moving cost, while a co-operative one is to only collect the fruit it prefers to maximize the social welfare of agents. In Section 4.4, we find that the two games satisfy the definition of SPD games in Section 2.3 by using policies with different cooperation degrees to play with each other.

## 4.2 Network Architecture and Parameter Settings

In both games, our network architectures for training the baseline policies follow standard AC networks, except that we allow both actor and critic to share the same underlying network to reduce the parameter space. For the underlying network, the first hidden layer convolves 32 filters of $8 \times 8$ with stride 4 with the input image and applies a rectifier non-linearity. The second hidden layer convolves 64 filters of $4 \times 4$ with stride 2, again followed by a rectifier nonlinearity. This is followed by a third convolutional layer that convolves 64 filters of $3 \times 3$ with stride 1 followed by a rectifier. For the

actor, on the basis of sharing network, the next layer includes 128 units with rectifier nonlinearity, and the final softmax layer has as many units as the number of actions. The critic is similar to the actor, but with only one scalar output.

The recurrent cooperation degree detection network is shown in Figure 1. The autoencoder and the detection network share the same underlying network. The first hidden layer convolves 10 filters of $3 \times 3$ with stride 2 with the input image and applies a rectifier nonlinearity. The second hidden layer is the same as the first one and the third hidden layer convolves 10 filters of $3 \times 3$ with stride 3 and applies a rectifier nonlinearity. The autoencoder is followed by a fourth hidden layer that deconvolves 10 filters of $3 \times 3$ with stride 3 and applies a sigmoid and its output shape is $21 \times 21 \times 10$. The next layer deconvolves 10 filters of $3 \times 3$ with stride 2 and applies a sigmoid and the output shape is $42 \times 42 \times 10$. The final layer deconvolves 10 filters of $3 \times 3$ with stride 2 and applies a sigmoid and the output shape is $84 \times 84 \times 3$. Cooperation degree detection network is followed by two LSTM layers of 256 units. The final layer is an output node.[1]

For the Apple-Pear game, each episode has at most 100 steps. The exploration rate is annealed linearly from 1 to 0.1 over the first 20000 steps. The weight parameters are updated by soft target updates [Lillicrap et al., 2015] every 4 steps to avoid the update fluctuation as follows:

$$\theta' \leftarrow 0.05\theta + 0.95\theta' \tag{16}$$

where $\theta$ is the parameter of the policy network and $\theta'$ is the

---

[1] The code and network architectures will be available soon: https://goo.gl/3VnFHj.

parameter of the target network. The learning rate is set to 0.0001 and memory is set to 25000. The batch size is 128. For the loss function in the cooperation degree detection network, we set $w_1$ and $w_2$ as 1 and 2 (see Equation 14). When agents play with different opponents online, we assign the number of states visited to the length $n$ of the state sequence which is the input of cooperation detection network. The co-operation degree changing factor $\alpha$ is set to 1.

The Fruit Gathering game uses the same detection network architecture. The actor-critic uses independent policy networks. Each episode has at most 100 steps during training. The exploration rate and the memory are the same as the Apple-Pear game. The weight parameters are updates the same as the Apple-Pear game:

$$\theta' \leftarrow 0.001\theta + 0.999\theta' \qquad (17)$$

When agents play with different opponents online, we set state sequence length $n$ as 50 and the changing factor $\alpha$ is set to 0.02.

### 4.3 Effect of Baseline Policy Generation

For the Apple-Pear game, the baseline policies are trained using the JAC scheme. The individual rewards of each agent increase gradually as its attitude increases and the other agent's attitude decreases. The results also indicate that an agent can learn a defecting policy when its attitude that represents its relative importance on overall reward increases and vice versa. We set the attitude in Equation (10) between 0.1 and 0.9 to train baseline policies. The reason that attitudes are not set 0 and 1 is that these settings will cause an agent with meaningless policy, e.g, the agent whose attitude equals 0 makes no influence to total reward in Equation (10) and hence will always avoid collecting fruit. This would, in turn, affect the learned policy quality of the other agent (i.e., lazy agent problem [Perolat *et al.*, 2017]). Figure 4 shows the average rewards of agents under policies trained with different weighted target rewards. We also evaluate the IAC scheme and similar results can be obtained and we omit it here.

For the Fruit Gathering game, the learning policies are synthesized based on IAC. IAC is more efficient for training baseline policies in this game, relative to JAC, since the rewards of collecting the apple for both agents are the same. On the other hand, if we adopt a JAC approach, it might lead to the consequence that the agent with a higher weight will collect all the apples, while the other agent will not collect any apples. Similar results as the Apple-Pear game can be observed here and we omit it here.

### 4.4 Effect of Policy Generation

For the Apple-Pear game, the baseline cooperation policies $\pi_1^c$ and $\pi_2^c$ are trained under the setting of $\{att_1 = att_2 = 0.5\}$. For $agent_1$, the baseline defection policy $\pi_1^d$ has $\{att_1 = 0.75, att_2 = 0.25\}$. For $agent_2$, the baseline defection policy $\pi_2^d$ has $\{att_1 = 0.25, att_2 = 0.75\}$. Then we generate the policy set of $agent_1$ $\Pi_1 = \{\pi = w_1 \times \pi_1^c + (1 - w_1) \times \pi_1^d | w_1 = 0.0, 0.1, \ldots, 1.0\}$ and the policy set of $agent_2$ $\Pi_2 = \{\pi = w_2 \times \pi_2^c + (1 - w_2) \times \pi_2^d | w_2 = 0.0, 0.1, \ldots, 1.0\}$. After that, two policies $\pi_1$ and $\pi_2$ are sampled from $\Pi_1$ and $\Pi_2$ and matched against each other in the games for 200,000

episodes. The average rewards are assigned to individual cells of different attitude pairs, in which $\pi_1$ correspond to policies with varying cooperation degrees $w_1$ for agent 1, and $\pi_2$, policies with varying cooperation degrees $w_2$ for agent 2 (see Figure 5 (a)). In Figure 5 (a), we can observe that when agents' cooperation degrees decrease, their rewards decrease. When both of their cooperation degrees increase, which means they are more cooperative, the sum of their rewards increase. Besides, given a fixed cooperation degree of $agent_1$, $agent_2$'s reward is increased as its cooperation degree decreases, and vice versa. Similar pattern can be observed for agent 1 as well. Therefore, it indicates that we can successfully synthesize policies with a continuous range of cooperation degrees. It also confirms that the Apple-Pear game can be seen as an SPD following the definition in Section 2.3.

For the Fruit Gathering game, we use policies with both attitudes equal to 0 and 0.5 as the baseline defection and cooperation policy respectively. Figure 5 (b) shows the results of their rewards, which is similar to the Apple-Pear game.

### 4.5 Effect of Cooperation Degree Detection

This section evaluates the detection power of the cooperation degree detection network. First, we train the cooperation degree detection network using datasets which include only the data labeled as full cooperation or full defection. The training data $(\langle s_0, s_1, s_2, \ldots, s_n \rangle, label)$ in the simulation is obtained based on baseline policies of $agent_2$. We set its label as 1 when $agent_2$ uses its baseline cooperation policy and 0 when $agent_2$ uses its baseline defection policy. Then we use this dataset to train the detection network. After training we evaluate the detection accuracy by applying it to detecting the cooperation degree of agent 2 when it uses policies with cooperation degrees from 0 to 1 and the degree interval is 0.1. The policy pair $(\pi_1, \pi_2)$ is sampled from $\Pi_1$ and $\Pi_2$ and matched against with each other for each episode. After the cooperation degree average value is stable, we view the output value as the cooperation degree of $agent_2$.

For the Apple-Pear game, we collect 10000 data for state sequence lengths $\{3, 4, 5, 6, 7, 8\}$. The cooperation detection results of $agent_2$ are shown in Figure 6 (a). We can see that the cooperation detection values can well approximate the true values with slight variance. Besides, the network can clearly detect the order of different cooperation degrees, which allows us to calculate the cooperation degree of $agent_2$ accurately. For the Fruit Gathering game, we collect 4000 data for each state sequence lengths $\{40, 50, 60, 70\}$, which includes 2000 data labeled as 1 and 2000 data labeled as 0. The network is trained in a similar way and the cooperation degree detection accuracy is high (see Figure 6 (b)). From Figure 6, We observe that the detection results are almost linear with the true values. Fig 6(c) shows the detection results of $agent_2$ when $agent_1$'s policy is fixed with cooperation degree 1. We can see that the true cooperation degree of $agent_2$ can be easily obtained by fitting a linear curve between the predicted and true values. Thus for each policy of $agent_1$ in $\Pi_1$, we can fit a linear function to evaluate the true cooperation degrees of $agent_2$. During practical online learning, when $agent_1$ uses policies of varying cooperation degrees to play with $agent_2$, it firstly chooses the function
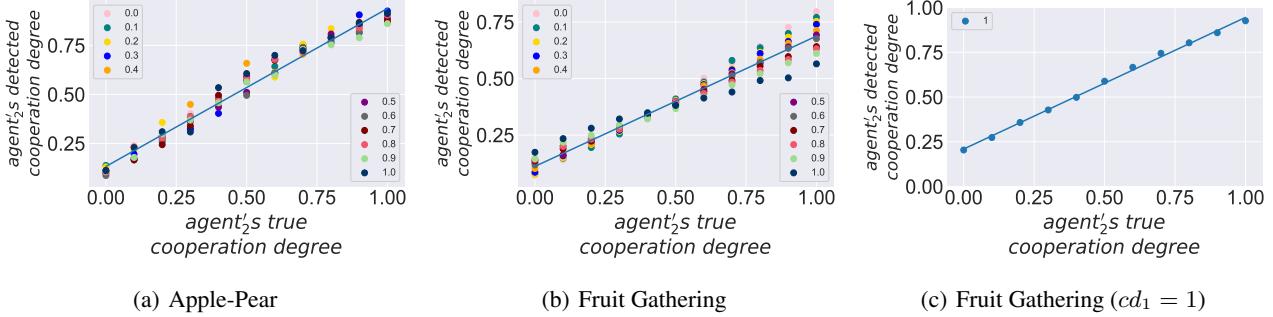
Figure 6: The detection results for $agent_2$ under different cooperation degrees of $agent_1$: (a) Apple-Pear game; (b) Fruit Gathering game; (c)$agent_1$'s cooperation degree as 1 in Fruit Gathering game

whose corresponding policy is closest to the used policy, and then computes the cooperation degree of $agent_2$.



Figure 7: Performance under self-play in the Apple-Pear game when both agents use our strategy.



Figure 8: Performance under self-play in the Fruit Gathering game when both agents use our strategy.

## 4.6 Performance under Self-Play

Next, we evaluate the learning performance of our approach under self-play. Since the initialized policies of agents can affect their behaviors in the game, we evaluate under all different initial conditions: a) $agent_1$ starts with cooperation policy and $agent_2$ starts with defecting policy; b) $agent_1$ starts with defecting policy and $agent_2$ starts with cooperation policy; c) both agents start with cooperation policies; d) both agents start with defecting policies. Agents converge to full cooperation for all four cases. We present the results for the last case, which is the most challenging one (see Figure 7 and Figure 8). When both agents start with a defection policy, it is more likely for them to model each other as defective and

thus both play defect thereafter. Our approach enables agents to successfully detect the cooperation tendency of their opponents from sequential actions and converge to cooperation at last. In the Apple-Pear game, agents converge efficiently within few episodes. The reason is that when agents are close to fruits they prefer, they will collect fruits no matter whether their policies are cooperative or not. Thus, an agent is more likely to be detected as being cooperative, which induces its opponent to change policies towards cooperation. In contrast, for the Fruit Gathering game, agents need a relatively longer time before converging to full cooperation. This is because in the Fruit Gathering game, the main feature of detecting the cooperation degree is beam emitting frequency. When one agent emits a beam continuously, they change to defect. Only when both agents collect fruits without emitting beam would lead to mutual cooperation.

## 4.7 Playing with Opponents with Changing Strategies

Now, we evaluate the performance against switching opponents, during a repeated interaction of $T$ episodes, the opponent changes its policy after a certain number of episodes and the learning agent does not know when the switches happen. Through this, we can evaluate the cooperation degree detection performance of our network, and verify whether our strategy can perform better than the fully cooperation or defection strategy from two aspects: 1) our approach seeks for mutual cooperation whenever possible; 2) our approach is robust against selfish exploitation. In the Apple-Pear game, we vary the value of $N_{change}$ in the range of $[50, 200]$ at the interval of 30, and similarly in the Fruit Gathering game we vary it in the range of $[100, 500]$ at the interval of 100. Only one set of results for each game are provided in Figure 9 and 10. The results of other values of $N_{change}$ are in the Appendix. Similar phenomenon can be observed for other values of $N_{change}$. From Figure 9 and 10, we observe that for both games, the average rewards of the learning agent ($agent_1$) are higher than its rewards using cooperation strategy $\pi^c$. And the social welfare (the sum of both agents' rewards) is higher than that using defection strategy $\pi^d$. This indicates that our approach can prevent agents from being exploited by defecting opponents and seek for cooperation against cooperative ones. By comparing the results for different values of $N_{change}$, we
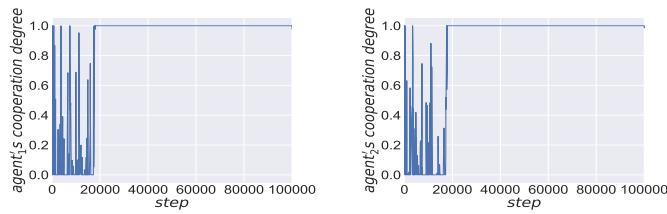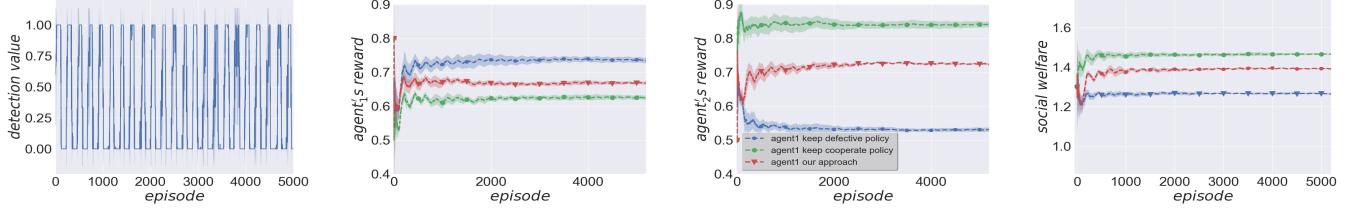
Figure 9: Apple-Pear game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 110 episodes. The average rewards of the $agent_1$ are higher when using our approach than using $\pi^c$, which means our approach can avoid being exploited by defective opponents. The social welfare is higher than using $\pi^d$, indicating that our approach can seek for cooperation against cooperative ones.
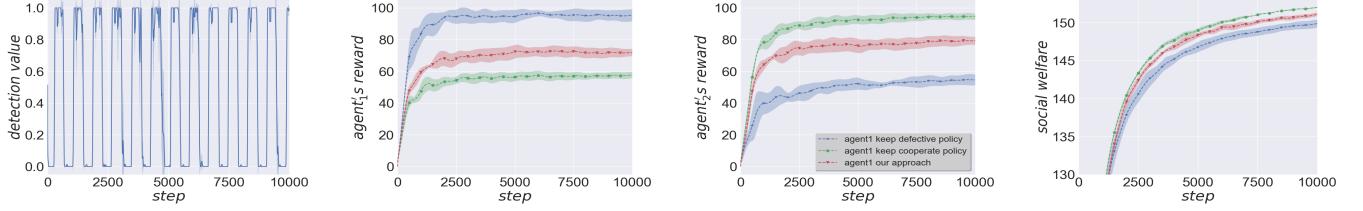


Figure 10: Fruit Gathering game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 300 steps. Similar phenomenon can be observed as in Figure 16.

find that the detection accuracy decreases when the opponent changes its policy quickly. Since the agent requires observing several episodes to detect the cooperation degree of the opponent, when the agent realizes its opponent changes the policy and adjusts its policy, the opponent may change the policy again. This problem becomes less severe when $T$ is comparatively large.

## 5 Conclusions

In this paper, we make the first step in investigating multiagent learning problem in large-scale PD games by leveraging the recent advance of deep reinforcement learning. A deep multiagent RL approach is proposed towards mutual cooperation in SPD games to support adaptive end-to-end learning. Empirical simulation shows that our agent can efficiently achieve mutual cooperation under self-play and also perform well against opponents with changing strategies. As the first step towards solving multiagent learning problem in large-scale environments, we believe there are many interesting questions remaining for future work. One worthwhile direction is how to generalize our approach to other classes of large-scale multiagent games. Generalized policy detection and reuse techniques should be proposed, e.g., by extending existing approaches in traditional reinforcement learning contexts [Hernandez-Leal *et al.*, 2016b; Hernandez-Leal *et al.*, 2016a; Hernandez-Leal and Kaisers, 2017].

## References

[Axelrod, 1984] Robert Axelrod. The evolution of cooperation, 1984.

[Banerjee and Sen, 2007] Dipyaman Banerjee and Sandip Sen. Reaching pareto-optimality in prisoner's dilemma using conditional joint action learning. *Autonomous Agents and Multi-Agent Systems*, 15(1):91–108, 2007.

[Busoniu *et al.*, 2008] Lucian Busoniu, Robert Babuska, and Bart De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, And Cybernetics-Part C: Applications and Reviews, 38 (2), 2008*, 2008.

[Crandall and Goodrich, 2005] Jacob W Crandall and Michael A Goodrich. Learning to teach and follow in repeated games. In *AAAI workshop on Multiagent Learning*, 2005.

[Crandall, 2012] Jacob W Crandall. Just add pepper: extending learning algorithms for repeated matrix games to repeated markov games. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*, pages 399–406. International Foundation for Autonomous Agents and Multiagent Systems, 2012.

[Damer and Gini, 2008] Steven Damer and Maria L Gini. Achieving cooperation in a minimally constrained environment. In *AAAI*, pages 57–62, 2008.

[Elidrisi *et al.*, 2014] Mohamed Elidrisi, Nicholas Johnson, Maria Gini, and Jacob Crandall. Fast adaptive learning in repeated stochastic games by game abstraction. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*, pages 1141–1148. International Foundation for Autonomous Agents and Multiagent Systems, 2014.

[Foerster *et al.*, 2017a] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. *arXiv preprint arXiv:1705.08926*, 2017.

[Foerster *et al.*, 2017b] Jakob Foerster, Nantas Nardelli, Gregory Farquhar, Philip Torr, Pushmeet Kohli, Shimon Whiteson, et al. Stabilising experience replay for deep multi-agent reinforcement learning. *arXiv preprint arXiv:1702.08887*, 2017.

[Hao and Leung, 2015] Jianye Hao and Ho-fung Leung. Introducing decision entrustment mechanism into repeated bilateral agent interactions to achieve social optimality. *Autonomous Agents and Multi-Agent Systems*, 29(4):658–682, 2015.

[He *et al.*, 2016] He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent modeling in deep reinforcement learning. In *International Conference on Machine Learning*, pages 1804–1813, 2016.

[Hernandez-Leal and Kaisers, 2017] Pablo Hernandez-Leal and Michael Kaisers. Towards a fast detection of opponents in repeated stochastic games. In *The Workshop on Transfer in Reinforcement Learning*, 2017.

[Hernandez-Leal *et al.*, 2016a] Pablo Hernandez-Leal, Benjamin Rosman, Matthew E Taylor, L Enrique Sucar, and Enrique Munoz de Cote. A bayesian approach for learning and tracking switching, non-stationary opponents. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, pages 1315–1316, 2016.

[Hernandez-Leal *et al.*, 2016b] Pablo Hernandez-Leal, Matthew E Taylor, Benjamin Rosman, L Enrique Sucar, and Enrique Munoz De Cote. Identifying and tracking switching, non-stationary opponents: a bayesian approach. In *Multiagent Interaction without Prior Coordination Workshop at AAAI*, 2016.

[Hu and Wellman, 2003] Junling Hu and Michael P Wellman. Nash q-learning for general-sum stochastic games. *Journal of machine learning research*, 4(Nov):1039–1069, 2003.

[Leibo *et al.*, 2017] Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. Multi-agent reinforcement learning in sequential social dilemmas. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, pages 464–473. International Foundation for Autonomous Agents and Multiagent Systems, 2017.

[Lillicrap *et al.*, 2015] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

[Littman, 1994] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Proceedings of the eleventh international conference on machine learning*, volume 157, pages 157–163, 1994.

[Lowe *et al.*, 2017] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *arXiv preprint arXiv:1706.02275*, 2017.

[Mathieu and Delahaye, 2015] Philippe Mathieu and Jean-Paul Delahaye. New winning strategies for the iterated prisoner's dilemma. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pages 1665–1666. International Foundation for Autonomous Agents and Multiagent Systems, 2015.

[Mnih *et al.*, 2013] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.

[Mnih *et al.*, 2015] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

[Nowak and Sigmund, 1993] Martin Nowak and Karl Sigmund. A strategy of win-stay, lose-shift that outperforms tit-for-tat in the prisoner's dilemma game. *Nature*, 364(6432):56–58, 1993.

[Perolat *et al.*, 2017] Julien Perolat, Joel Z Leibo, Vinicius Zambaldi, Charles Beattie, Karl Tuyls, and Thore Graepel. A multi-agent reinforcement learning model of common-pool resource appropriation. *arXiv preprint arXiv:1707.06600*, 2017.

[Schulman *et al.*, 2015] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.

[Stone and Veloso, 2000] Peter Stone and Manuela Veloso. Multiagent systems: A survey from a machine learning perspective. *Autonomous Robots*, 8(3):345–383, 2000.

[Sunehag *et al.*, 2017] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.

[Sutton *et al.*, 2000] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063, 2000.

[Wang *et al.*, 2016] Ziyu Wang, Victor Bapst, Nicolas Heess, Volodymyr Mnih, Remi Munos, Koray Kavukcuoglu, and Nando de Freitas. Sample efficient actor-critic with experience replay. *arXiv preprint arXiv:1611.01224*, 2016.

[Williams, 1992] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

# A  Playing with Opponents with Changing Strategies
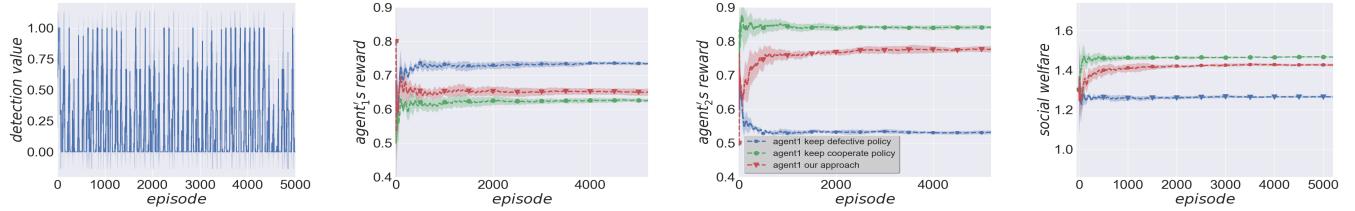
## A.1  The Apple-Pear Game



Figure 11: Apple-Pear game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 50 episodes.
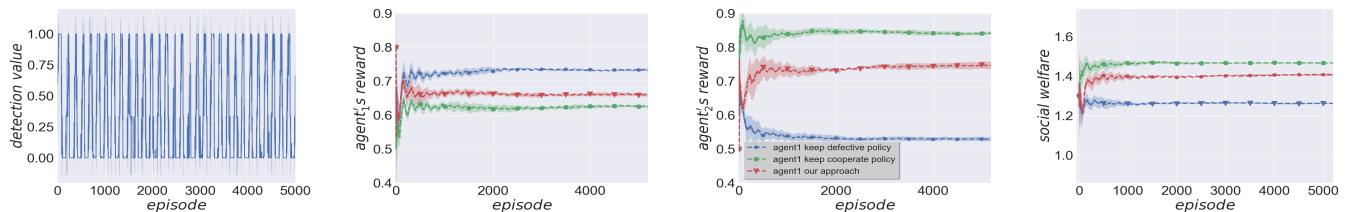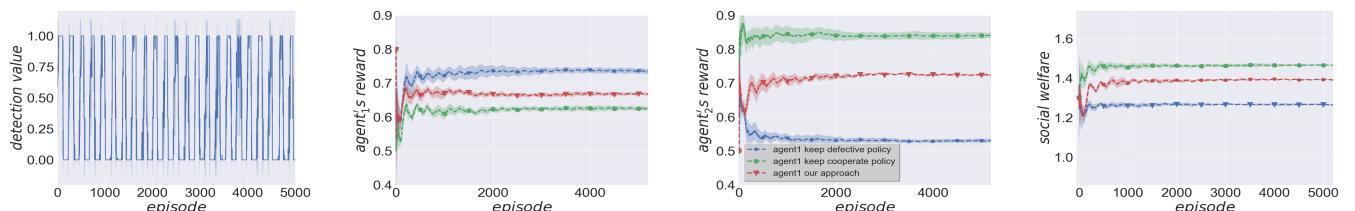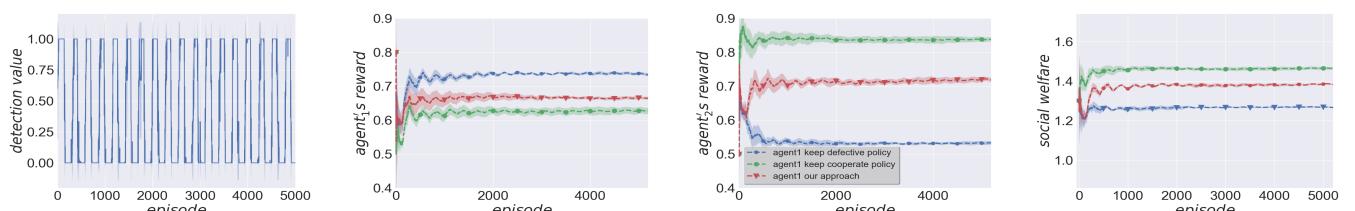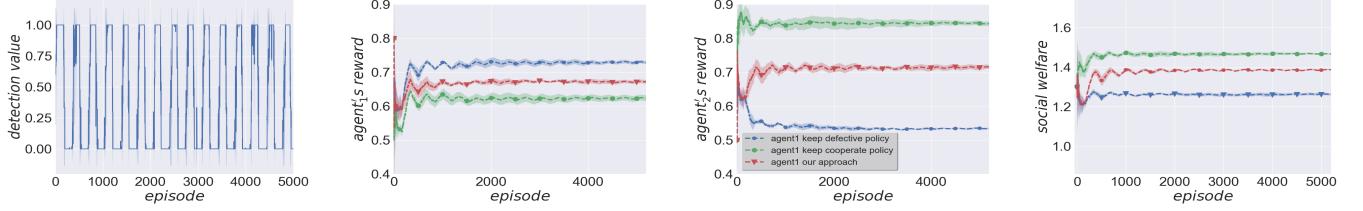


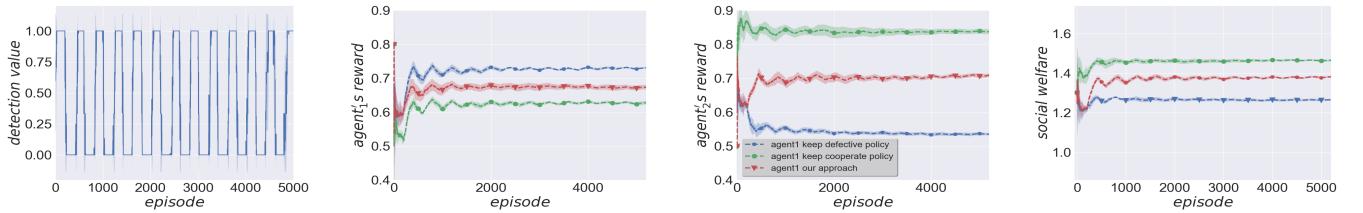Figure 12: Apple-Pear game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 80 episodes.



Figure 13: Apple-Pear game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 110 episodes.



Figure 14: Apple-Pear game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 140 episodes.

Figure 15: Apple-Pear game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 170 episodes.



Figure 16: Apple-Pear game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 200 episodes.
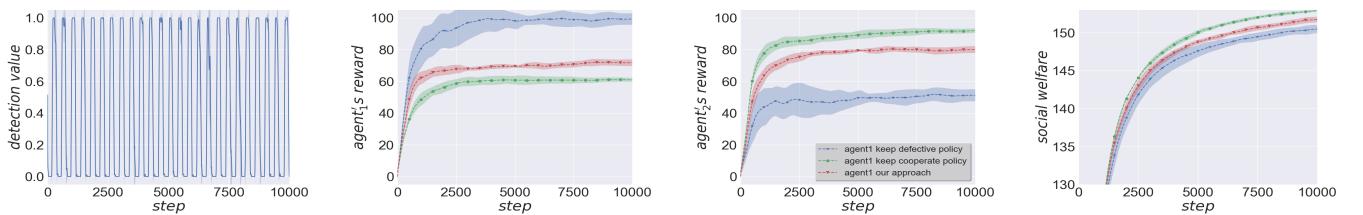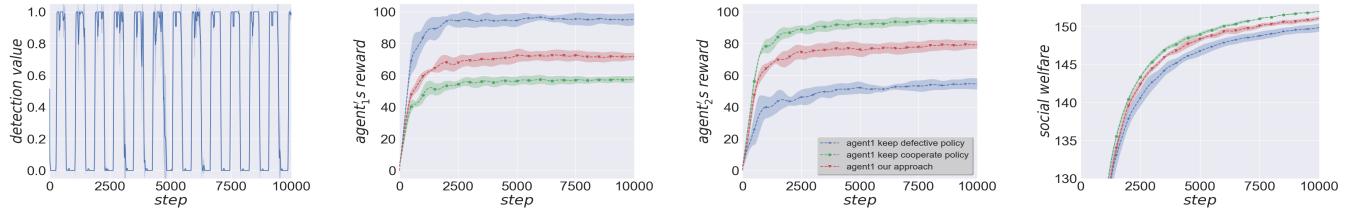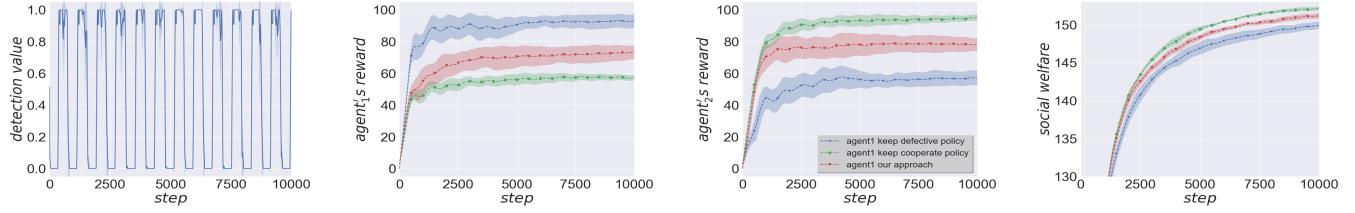
## A.2 The Gathering Game



Figure 17: Fruit Gathering game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 100 steps.



Figure 18: Fruit Gathering game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 200 steps.
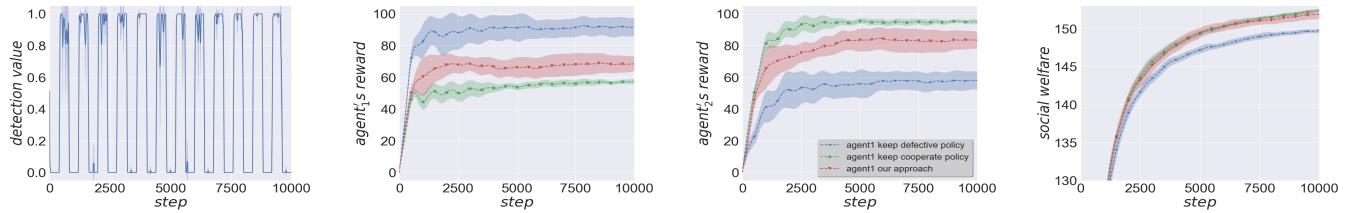
Figure 19: Fruit Gathering game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 300 steps.



Figure 20: Fruit Gathering game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 400 steps.



Figure 21: Fruit Gathering game: $agent_2$'s policy varies between $\pi^c$ and $\pi^d$ every 500 steps.