

Inducing Cooperative behaviour in Sequential-Social dilemmas through Multi-Agent Reinforcement Learning using Status-Quo Loss

Pinkesh Badjatiya¹ Mausoom Sarkar¹ Abhishek Sinha¹ Siddharth Singh² Nikaash Puri¹
Jayakumar Subramanian¹ Balaji Krishnamurthy¹

Abstract

In social dilemma situations, individual rationality leads to sub-optimal group outcomes. Several human engagements can be modeled as a sequential (multi-step) social dilemmas. However, in contrast to humans, Deep Reinforcement Learning agents trained to optimize individual rewards in sequential social dilemmas converge to selfish, mutually harmful behavior. We introduce a status-quo loss (*SQLoss*) that encourages an agent to stick to the status-quo, rather than repeatedly changing its policy. We show how agents trained with *SQLoss* evolve cooperative behavior in several social dilemma matrix games. To work with social dilemma games that have visual input, we propose *GameDistill*. *GameDistill* uses self-supervision and clustering to automatically extract cooperative and selfish policies from a social dilemma game. We combine *GameDistill* and *SQLoss* to show how agents evolve socially desirable cooperative behavior in the Coin Game.

1. Introduction

Consider a sequential social dilemma, where individually rational behavior leads to outcomes that are sub-optimal for each individual in the group. (Hardin, 1968; Ostrom, 1990; Ostrom et al., 1999; Dietz et al., 2003). Current state-of-the-art Multi-Agent Deep Reinforcement Learning (MARL) methods that train agents independently can lead to agents that fail to cooperate reliably, even in simple social dilemma settings. This failure to cooperate results in sub-optimal individual and group outcomes (Foerster et al. (2018); Lerer & Peysakhovich (2017), Section 2.2).

To illustrate why it is challenging to evolve cooperation in such dilemmas, we consider the Coin Game ((Foerster et al.,

¹Media and Data Science Research Lab, Adobe ²IIT Kharagpur, India. Correspondence to: Pinkesh Badjatiya <pbad-jati@adobe.com>.

2018), Figure 1). Each agent can play either selfishly (pick all coins) or cooperatively (pick only coins of its color). Regardless of the behavior of the other agent, the individually rational choice for an agent is to play selfishly, either to minimize losses (avoid being exploited) or to maximize gains (exploit the other agent). However, when both agents behave rationally, they try to pick all coins and achieve an average long term reward of -0.5 . In contrast, if both play cooperatively, then the average long term reward for each agent is 0.5 . Therefore, when agents cooperate, they are both better off. Training Deep RL agents independently in the Coin Game using state-of-the-art methods leads to mutually harmful selfish behavior (Section 2.2).

In this paper, we present a novel MARL algorithm that allows independently learning Deep RL agents to converge to **individually and socially desirable cooperative behavior** in such social dilemma situations. Our key contributions can be summarised as:

1. We introduce a **Status-Quo** loss (*SQLoss*, Section 2.3) and an associated policy gradient-based algorithm to evolve optimal behavior for agents that can act in either a cooperative or a selfish manner, by choosing between a cooperative and a selfish policy. We empirically demonstrate that agents trained with the *SQLoss* evolve optimal behavior in several social dilemma iterated matrix games (Section 4).
2. We propose *GameDistill* (Section 2.4), an algorithm that reduces a social dilemma game with visual observations to an iterated matrix game by extracting policies that implement cooperative and selfish behavior. We empirically demonstrate that *GameDistill* extracts cooperative and selfish policies for the Coin Game (Section 4.2).
3. We demonstrate that when agents run *GameDistill* followed by MARL game-play using *SQLoss*, they converge to individually as well as socially desirable cooperative behavior in a social dilemma game with visual observations (Section 4.2).

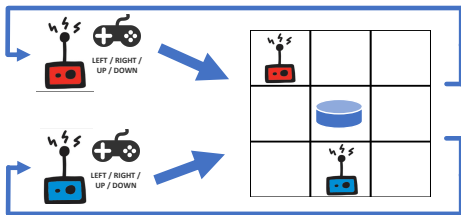


Figure 1: Two agents (Red and Blue) playing the Coin Game. The agents, along with a Blue or Red coin, appear at random positions in a 3x3 grid. An agent observes the complete 3x3 grid as input and can move either left, right, up, or down. When an agent moves into a cell with a coin, it picks the coin, and a new instance of the game begins. If the Red agent picks the Red coin, it gets a reward of +1 and the Blue agent gets no reward. If the Red agent picks the Blue coin, it gets a reward of +1, and the Blue agent gets a reward of -2 . The blue agent’s reward structure is symmetric to that of the red agent.

cooperative behavior in social dilemmas has been studied by researchers through human studies and simulation models (Fudenberg & Maskin, 1986; Green & Porter, 1984; Fudenberg et al., 1994; Kamada & Kominers, 2010; Abreu et al., 1990). A large body of work has looked at the mechanism of evolution of cooperation through reciprocal behaviour and indirect reciprocity (Trivers, 1971; Axelrod, 1984; Nowak & Sigmund, 1992; 1993; 1998), through variants of reinforcement using aspiration (Macy & Flache, 2002), attitude (Damer & Gini, 2008) or multi-agent reinforcement learning (Sandholm & Crites, 1996; Wunder et al., 2010), and under specific conditions (Banerjee & Sen, 2007) using different learning rates (de Cote et al., 2006) similar to WoLF (WOL, 2002) as well as using embedded emotion (Yu et al., 2015), social networks (Ohtsuki et al., 2006; Santos & Pacheco, 2006).

However, these approaches do not directly apply to Deep RL agents (Leibo et al., 2017). Recent work in this direction (Kleiman-Weiner et al., 2016; Julien et al., 2017; Peysakhovich & Lerer, 2018) focuses on letting agents learn strategies in multi-agent settings through interactions with other agents. Leibo et al. (2017) define the problem of social dilemmas in the Deep RL framework and analyze the outcomes of a fruit-gathering game (Julien et al., 2017). They vary the abundance of resources and the cost of conflict in the fruit environment to generate degrees of cooperation between agents. Hughes et al. (2018) define an intrinsic reward (inequality aversion) that attempts to reduce the difference in obtained rewards between agents. The agents are designed to have an aversion to both advantageous (guilt) and disadvantageous (unfairness) reward allocation. This handcrafting of loss with mutual fairness evolves cooperation, but it leaves the agent vulnerable to exploitation. LOLA (Foerster

et al., 2018) uses opponent awareness to achieve high levels of cooperation in the Coin Game and the Iterated Prisoner’s Dilemma game. However, the LOLA agent assumes access to the other agent’s policy parameters and gradients. This level of access is analogous to getting complete access to the other agent’s private information and therefore devising a strategy with full knowledge of how they are going to play. Wang et al. (2019) propose an evolutionary Deep RL setup to evolve cooperation. They define an intrinsic reward that is based on features generated from the agent’s past and future rewards, and this reward is shared with other agents. They use evolution to maximize the sum of rewards among the agents and thus evolve cooperative behavior. However, sharing rewards in this indirect way enforces cooperation rather than evolving it through independently learning agents.

In contrast, we introduce a Status-Quo ($SQLoss$) that evolves cooperation between agents without sharing rewards, gradients, or using a communication channel. The $SQLoss$ encourages an agent to imagine the consequences of sticking to the status-quo. This imagined stickiness ensures that an agent gets a better estimate of a cooperative or selfish policy. Without $SQLoss$, agents repeatedly switch policies (from cooperative to selfish), obtain short-term rewards (through exploitation), and, therefore incorrectly learn that a selfish strategy gives higher rewards in the long-term.

To work with social dilemma games that have visual observations, we introduce *GameDistill*. *GameDistill* uses self-supervision and clustering to automatically extract a cooperative and a selfish policy from a social dilemma game. The input to *GameDistill* is a collection of state sequences derived from game-play between two randomly initialized agents. Each state sequence represents a collection of states and actions (of both agents), leading up to a reward in the environment. *GameDistill* uses this collection of state sequences to learn two oracles. One oracle represents a cooperative policy, and the other oracle represents a selfish policy. Given a state, an oracle returns an action according to the specific policy. It is important to note that each agent independently runs *GameDistill* to extract oracles. (For instance, Figure 8 (Appendix A) illustrates the cooperation and defection oracles extracted by the Red agent using *GameDistill* in the Coin Game.) Figure 2 shows the high-level architecture of our approach.

1. For a social dilemma game with visual observations, each RL agent runs *GameDistill* to learn oracles that implement cooperative and selfish behavior.
2. We train agents (with $SQLoss$) to play the game such that at any step, an agent can either take the action suggested by the cooperation oracle or the selfish oracle.

We empirically demonstrate in Section 4 that our approach

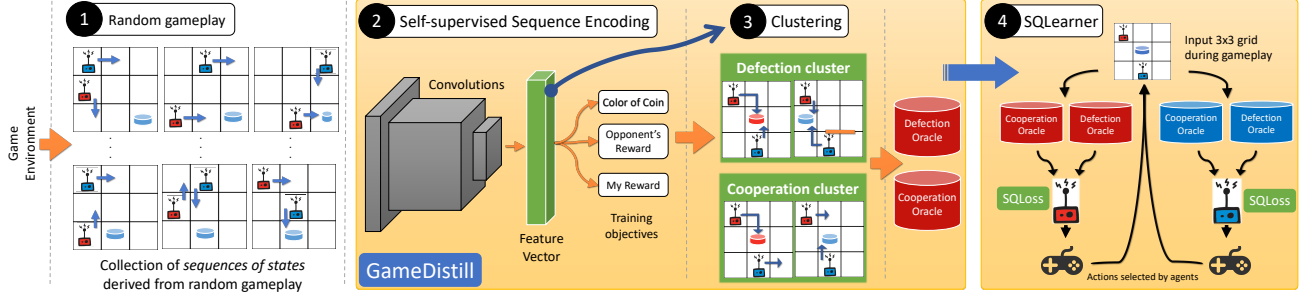


Figure 2: High-level architecture of our approach. Each agent runs *GameDistill* by performing (1), (2), (3) individually before playing the social dilemma game. This creates two oracles per agent. During game-play(4), each agent (enhanced with *SQLoss*) takes either the action suggested by the cooperation oracle or the action suggested by the defection oracle.

evolves cooperative behavior between independently trained agents.

2. Approach

2.1. Social Dilemmas modeled as Iterated Matrix Games

We adopt the definitions in Foerster et al. (2018). We model social dilemmas as general-sum Markov (simultaneous move) games. A multi-agent Markov game is specified by $G = \langle S, A, U, P, r, n, \gamma \rangle$. S denotes the state space of the game. n denotes the number of agents playing the game. At each step of the game, each agent $a \in A$, selects an action $u^a \in U$. \vec{u} denotes the joint action vector that represents the simultaneous actions of all agents. The joint action \vec{u} changes the state of the game from s to s' according to the state transition function $P(s'|\vec{u}, s) : S \times \mathbf{U} \times S \rightarrow [0, 1]$. At the end of each step, each agent a gets a reward according to the reward function $r^a(s, \vec{u}) : S \times \mathbf{U} \rightarrow \mathbb{R}$. The reward obtained by an agent at each step is a function of the actions played by all agents. For an agent, a , the discounted future return from time t is defined as $R_t^a = \sum_{l=0}^{\infty} \gamma^l r_{t+l}^a$, where $\gamma \in [0, 1)$ is the discount factor. Each agent independently attempts to maximize its expected total discounted reward.

Matrix games are the special case of two-player perfectly observable Markov games (Foerster et al., 2018). Table 1 shows examples of matrix games that represent social dilemmas. Consider the Prisoner's Dilemma matrix game in Table 1a. Each agent can either cooperate (C) or defect (D). For an agent, playing D is the rational choice, regardless of whether the other agent plays either C or D . Therefore, if both agents play rationally, they each receive a reward of -2 . However, if each agent plays C , then it will obtain a reward of -1 . This fact that individually rational behavior leads to a sub-optimal group (and individual) outcome highlights the dilemma.

In Infinitely Iterated Matrix Games, agents repeatedly play a particular matrix game against each other. In each iteration

of the game, each agent has access to the actions played by both agents in the previous iteration. Therefore, the state input to an RL agent consists of the actions of both agents in the previous iteration of the game. We adopt this state formulation to remain consistent with Foerster et al. (2018). The infinitely iterated variations of the matrix games in Table 1 represent sequential social dilemmas. For ease of representation, we refer to infinitely iterated matrix games as iterated matrix games in subsequent sections.

2.2. Learning Policies in Iterated Matrix Games: The Selfish Learner

The standard method to model agents in iterated matrix games is to model each agent as a Deep RL agent that independently attempts to maximize its expected total discounted reward. Several approaches to model agents in this way use policy gradient-based methods (Sutton et al., 2000; Williams, 1992)). Policy gradient methods update an agent's policy, parameterized by θ^a , by performing gradient ascent on the expected total discounted reward $\mathbb{E}[R_0^a]$. Formally, let θ^a denote the parameterized version of an agent's policy π^a and V_{θ^1, θ^2}^a denote the total expected discounted reward for agent a . Here, V^a is a function of the policy parameters (θ^1, θ^2) of both agents. In the i^{th} iteration of the game, each agent updates θ_i^a to θ_{i+1}^a , such that it maximizes its total expected discounted reward. θ_{i+1}^a is computed as follows:

$$\begin{aligned} \theta_{i+1}^1 &= \operatorname{argmax}_{\theta^1} V^1(\theta^1, \theta_i^2) \\ \theta_{i+1}^2 &= \operatorname{argmax}_{\theta^2} V^2(\theta_i^1, \theta^2) \end{aligned} \quad (1)$$

For agents trained using reinforcement learning, the gradient ascent rule to update θ_{i+1}^1 is:

$$\begin{aligned} f_{nl}^1 &= \nabla_{\theta_i^1} V^1(\theta_i^1, \theta_i^2) \cdot \delta \\ \theta_{i+1}^1 &= \theta_i^1 + f_{nl}^1(\theta_i^1, \theta_i^2) \end{aligned} \quad (2)$$

where δ is the step size of the updates.

In the Iterated Prisoner's Dilemma (IPD) game, agents trained with the policy gradient update method converge

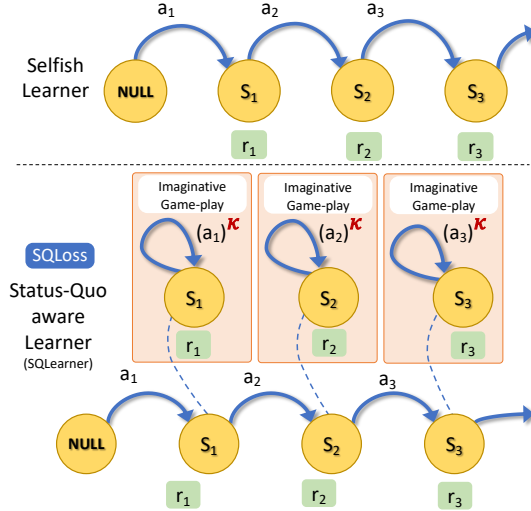


Figure 3: Intuition behind *SQLearner*. At each step, the *SQLoss* encourages an agent to imagine the consequences of sticking to the status-quo by imagining an episode where the status-quo is repeated for κ steps. Section 2.3 describes *SQLoss* in more detail.

to a sub-optimal mutual defection equilibrium (Figure 4, (Lerer & Peysakhovich, 2017)). This sub-optimal equilibrium attained by Selfish Learners motivates us to explore alternative methods that could lead to a desirable cooperative equilibrium. We denote the agent trained using policy gradient updates as a Selfish Learner (*SL*).

2.3. Learning Policies in Iterated Matrix Games: The Status-Quo Aware Learner (*SQLoss*)

2.3.1. *SQLoss*: INTUITION

Why do independent, selfish learners converge to mutually harmful behavior in the IPD? To understand this, consider the payoff matrix for a single iteration of the IPD in Table 1a. In each iteration, an agent can play either *C* or *D*. Mutual defection (*D, D*) is worse for each agent than mutual cooperation (*C, C*). However, one-sided exploitation (*D, C*) is better than mutual cooperation for the exploiter and far worse for the exploited. Therefore, as long as an agent perceives the possibility of exploitation (*(D, C)* or *(C, D)*), it is drawn to defect, both to maximize reward (through exploitation) and minimize loss (through being exploited). To increase the likelihood of cooperation, it is important to reduce instances of exploitation between agents. We posit that, if agents mostly only either mutually cooperate (*C, C*) or mutually defect (*D, D*), then they will learn to prefer *C* and achieve a socially desirable cooperative equilibrium.

Motivated by this idea, we introduce a status-quo loss (*SQLoss*) for each agent derived from the idea of imaginary game-play, as depicted in Figure 3. Intuitively, the

loss encourages an agent to imagine an episode where the status-quo (current situation) is repeated for a number of steps. This imagined episode causes the exploited agent (in *(D, C)*) to perceive a continued risk of exploitation and, therefore, quickly move to *(D, D)*. Hence, for an agent, the short-term gain from exploitation (*(D, C)*) is overcome by the long-term loss from mutual exploitation (*(D, D)*). Therefore, agents move towards either mutual cooperation (*(C, C)*) or mutual defection (*(D, D)*). With exploitation (and subsequently, the fear of being exploited) out of the picture, agents move towards mutual cooperation. Figure 3 shows the idea behind *SQLoss*.

2.3.2. *SQLoss*: FORMULATION

We describe below the formulation of *SQLoss* with respect to agent 1. The formulation for agent 2 is identical to that of agent 1. Let $\tau_a = (s_0, u_0^1, u_0^2, r_0^1, \dots, s_T, u_T^1, u_T^2, r_T^1)$ denote the collection of an agent’s experiences after T time steps. Let $R_t^1(\tau_1)$ denote the discounted future return for agent 1 starting at s_t in actual game-play. Let $\hat{\tau}_1$ denote the collection of an agent’s **imagined** experiences. For a state s_t ($t \in [0, T)$), an agent imagines an episode by starting at s_t and repeating u_{t-1}^1, u_{t-1}^2 for κ_t steps. This is equivalent to imagining a κ_t step repetition of already played actions. We sample κ_t from a Discrete Uniform distribution $\mathbb{U}\{1, z\}$ where z is a hyper-parameter ≥ 1 . To simplify notation, let $\phi_t(s_t, \kappa_t)$ denote the ordered set of state, actions, and rewards starting at time t and repeated κ_t times for imagined game-play. Let $\hat{R}_t^1(\hat{\tau}_1)$ denote the discounted future return starting at s_t in imagined status-quo game-play.

$$\phi_t(s_t, \kappa_t) = [(s_t, u_{t-1}^1, u_{t-1}^2, r_{t-1}^1)_0, \dots, (s_t, u_{t-1}^1, u_{t-1}^2, r_{t-1}^1)_{\kappa_t-1}] \quad (3)$$

$$\hat{\tau}_1 = (\phi_t(s_t, \kappa_t), (s_{t+1}, u_{t+1}^1, u_{t+1}^2, r_{t+1}^1)_{\kappa_t+1}, \dots, (s_T, u_T^1, u_T^2, r_T^1)_{T+\kappa_t-t}) \quad (4)$$

$$\begin{aligned} R_t^1(\tau_1) &= \sum_{l=t}^T \gamma^{l-t} r_l^1 \\ \hat{R}_t^1(\hat{\tau}_1) &= \left(\frac{1-\gamma^\kappa}{1-\gamma} \right) r_{t-1}^1 + \gamma^\kappa R_t^1(\tau_1) \\ &= \left(\frac{1-\gamma^\kappa}{1-\gamma} \right) r_{t-1}^1 + \gamma^\kappa \sum_{l=t}^T \gamma^{l-t} r_l^1 \end{aligned} \quad (5)$$

V_{θ^1, θ^2}^1 and $\hat{V}_{\theta^1, \theta^2}^1$ are approximated by $\mathbb{E}[R_0^1(\tau_1)]$ and $\mathbb{E}[\hat{R}_0^1(\hat{\tau}_1)]$ respectively. These V values are the expected rewards conditioned on both agents’ policies (π^1, π^2). For

agent 1, the regular gradients and the Status-Quo gradients, $\nabla_{\theta^1} \mathbb{E}[R_0^1(\tau_1)]$ and $\nabla_{\theta^1} \mathbb{E}[\hat{R}_0^1(\hat{\tau}_1)]$, can be derived from the policy gradient formulation as:

$$\begin{aligned} \nabla_{\theta^1} \mathbb{E}[R_0^1(\tau_1)] &= \mathbb{E}[R_0^1(\tau_1) \nabla_{\theta^1} \log \pi^1(\tau_1)] \\ &= \mathbb{E} \left[\sum_{t=0}^T \nabla_{\theta^1} \log \pi^1(u_t^1 | s_t) \cdot \sum_{l=t}^T \gamma^l r_l^1 \right] \\ &= \mathbb{E} \left[\sum_{t=0}^T \nabla_{\theta^1} \log \pi^1(u_t^1 | s_t) \gamma^t (R_t^1(\tau_1) - b(s_t)) \right] \end{aligned} \quad (6)$$

$$\begin{aligned} \nabla_{\theta^1} \mathbb{E}[\hat{R}_0^1(\hat{\tau}_1)] &= \mathbb{E} [\hat{R}_0^1(\hat{\tau}_1) \nabla_{\theta^1} \log \pi^1(\hat{\tau}_1)] \\ &= \mathbb{E} \left[\sum_{t=0}^T \nabla_{\theta^1} \log \pi^1(u_{t-1}^1 | s_t) \times \right. \\ &\quad \left. \left(\sum_{l=t}^{t+\kappa} \gamma^l r_{t-1}^1 + \sum_{l=t}^T \gamma^{l+\kappa} r_l^1 \right) \right] \\ &= \mathbb{E} \left[\sum_{t=0}^T \nabla_{\theta^1} \log \pi^1(u_{t-1}^1 | s_t) \times \right. \\ &\quad \left. \left(\left(\frac{1-\gamma^\kappa}{1-\gamma} \right) \gamma^t r_{t-1}^1 + \gamma^\kappa \sum_{l=t}^T \gamma^l r_l^1 \right) \right] \\ &= \mathbb{E} \left[\sum_{t=0}^T \nabla_{\theta^1} \log \pi^1(u_{t-1}^1 | s_t) \gamma^t (\hat{R}_t^1(\hat{\tau}_1) - b(s_t)) \right] \end{aligned} \quad (7)$$

where $b(s_t)$ is a baseline for variance reduction.

Then the update rule $f_{sql,pg}^1$ for the policy gradient-based Status-Quo Learner (SQL-PG) is,

$$f_{sql,pg}^1 = (\alpha \cdot \nabla_{\theta^1} \mathbb{E}[R_0^1(\tau_1)] + \beta \cdot \nabla_{\theta^1} \mathbb{E}[\hat{R}_0^1(\hat{\tau}_1)]) \cdot \delta \quad (8)$$

where α and β denote the loss scaling factor for the reinforce and the imaginative game-play, respectively.

2.4. *GameDistill*: Moving Beyond Iterated Matrix Games

In the previous sections, we have focused on evolving cooperative behavior in the iterated matrix game formulation of sequential social dilemmas. In the iterated matrix game formulation, an agent is only allowed to either cooperate or defect in each iteration. However, in a social dilemma game with visual observations, it is not clear what set of low-level actions constitute cooperative or selfish behavior. Therefore, to work on social dilemmas with visual observations, we propose *GameDistill*, an approach that automatically learns a cooperation and a defection policy by analyzing the behavior of randomly initialized agents. *GameDistill*

learns these policies in the form of cooperation and defection oracles. Given a state, the cooperation oracle suggests an action that represents cooperative behavior. Similarly, the defection oracle suggests an action that represents selfish behavior (Figure 8 (Appendix A)). When RL agents play the social dilemma game, each agent independently runs *GameDistill* before playing the game. Once both agents have run *GameDistill*, they consult either of the two extracted oracles in every step of the game. Therefore, in each step, an agent either takes the action recommended by the cooperation oracle or the action recommended by the defection oracle. In this way, we reduce the visual input game to an iterated matrix game and subsequently apply *SQLoss* to evolve cooperative behavior. *GameDistill* (see Figure 2) works as follows.

1. We initialize RL agents with random weights and play them against each other in the game. In these random game-play episodes, whenever an agent receives a reward, we store the sequence of the last three states up to the current state.
2. This collection of state sequences is used to train the *GameDistill* network. The *GameDistill* network takes as input a sequence of states and predicts the rewards of both agents as well as environment parameters that depend on the game. For instance, in the Coin Game, the *GameDistill* network predicts the rewards of both agents and the color of the picked coin.
3. Training the *GameDistill* network leads to the emergence of feature embeddings for the various state sequences. Subsequently, clustering these embeddings using Agglomerative Clustering (with number of clusters=2) leads to the development of cooperative and defection clusters. One of the learned clusters contains state sequences that represent cooperative behavior, and the other cluster contains state sequences that represent defection. For instance, in the Coin Game, a point in the cooperation cluster contains a sequence of states where an agent picks a coin of its color.
4. To train the cooperation and defection oracle networks, we use the collection of state sequences in each cluster. For each sequence of states in a cluster, we train the oracle network to predict the next action, given the current state. For instance, Figure 8 (Appendix A) shows the cooperation and defection oracles extracted by the Red agent using *GameDistill* in the Coin Game.

Section 3.3 describes the architectural choices of each component of *GameDistill*.

	<i>C</i>	<i>D</i>		<i>H</i>	<i>T</i>		<i>C</i>	<i>D</i>
<i>C</i>	(-1, -1)	(-3, 0)	<i>H</i>	(+1, -1)	(-1, +1)	<i>C</i>	(0, 0)	(-4, -1)
<i>D</i>	(0, -3)	(-2, -2)	<i>T</i>	(-1, +1)	(+1, -1)	<i>D</i>	(-1, -4)	(-3, -3)

(a) Prisoners' Dilemma (PD) (b) Matching Pennies (MP) (c) Stag Hunt (SH)

Table 1: Payoff matrices for the different games used in our experiments. (X, Y) in a cell represents a reward of X to the row and Y to the column player. C , D , H , and T denote the actions for the row and column players. In the iterated versions of these games, agents play against each other over several iterations. In each iteration, an agent takes an action and receives a reward based on the actions of both agents. Each matrix represents a different kind of social dilemma.

3. Experimental Setup

In order to compare our results to previous work (Foerster et al., 2018), we use the Normalized Discounted Reward ($NDR = (1 - \gamma) \sum_{t=0}^T \gamma^t r_t$). A higher NDR implies that an agent obtains a higher reward in the environment. We compare our approach (Status-Quo Aware Learner: *SQ Learner*) to Learning with Opponent-Learning Awareness (Lola-PG) (Foerster et al., 2018) and the Selfish Learner (SL, Section 2.2). For all experiments, we perform 20 runs and report average *NDR*, along with variance across runs. The bold line in all the figures is the mean, and the shaded region is the one standard deviation region over the mean. All of our code is available at Code (2019).

3.1. Social Dilemma Games

For our experiments with social dilemma matrix games, we use the (Iterated Prisoners Dilemma (IPD) (Luce & Raiffa, 1989), Iterated Matching Pennies (IMP) (Lee & Louis, 1967), and the Iterated Stag Hunt (ISH) (Foerster et al., 2018)). Table 1 shows the payoff matrix for a single iteration of each game. In iterated matrix games, at each iteration, agents take an action according to a policy and receive the rewards in Table 1. To simulate an infinitely iterated game, we let agents play 200 iterations of the game against each other, and do not provide an agent with any information about the number of remaining iterations (Foerster et al., 2018). In an iteration, the state for an agent is the actions played by both agents in the previous iteration. Each matrix game in Table 1 represents a different dilemma.

In the Prisoner's Dilemma, the rational policy for each agent is to defect, regardless of the policy of the other agent. However, when each agent plays rationally, each is worse off. In Matching Pennies, if an agent plays predictably, it is prone to exploitation by the other agent. Therefore, the optimal policy is to randomize between H and T , obtaining an average NDR of 0. The Stag Hunt game represents a coordination dilemma. In the game, given that the other agent will cooperate, an agent's optimal action is to cooperate as well. However, at each step, each agent has an attractive alternative, that of defecting and obtaining a guaranteed reward of 3. Therefore, the promise of a safer alternative and the

fear that the other agent might select the safer choice could drive an agent to also select the safer alternative, thereby sacrificing the higher reward of mutual cooperation.

For our experiments on a social dilemma with visual observations, we use the Coin Game (Figure 1) (Foerster et al., 2018). The rational policy for an agent is to defect and try to pick all coins, regardless of the policy of the other agent. However, when both agents defect, both are worse off.

3.2. SQLoss

For our experiments with the Selfish and Status-Quo Aware Learner (*SQ Learner*), we use policy gradient-based learning where we train an agent with the Actor-Critic method (Sutton & Barto, 2011). Each agent is parameterized with a policy actor and critic for variance reduction in policy updates. During training, we use gradient descent with step size, $\delta = 0.005$ for the actor and $\delta = 1$ for the critic. We use a batch size of 4000 for Lola-PG (Foerster et al., 2018) and 200 for *SQ Learner* for roll-outs. We use an episode length of 200 for all iterated matrix games. We use a discount rate (γ) of 0.96 for the Iterated Prisoners' Dilemma, Iterated Stag Hunt, and Coin Game. For the Iterated Matching Pennies, we use $\gamma = 0.9$. The high value of γ allows for long time horizons, thereby incentivizing long-term reward. Each agent randomly samples κ from $\mathbb{U} \in (1, z)$ ($z = 10$, discussed in Appendix C) at each step.

3.3. GameDistill

GameDistill consists of two components. First, the state sequence encoder (Step 2, Section 2.4) that takes as input a sequence of states and outputs a feature representation. We encode each state in the sequence using a series of standard Convolution layers with kernel-size 3. We then use a fully-connected layer with 100 neurons that outputs a dense representation of the sequence of states. The color of picked coin, agent reward, and opponent reward branches consist of a series of dense layers with linear activation. We use linear activation so that we can cluster the feature vectors (embeddings) using a linear clustering algorithm, such as Agglomerative Clustering. We obtain similar results when we use the K-means clustering algorithm. We use the

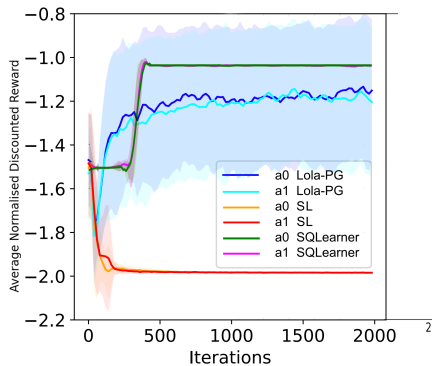


Figure 4: Average NDR values for different learners in the IPD game. *SQLearner* agents obtain a near-optimal NDR value (-1) for this game. In contrast to other methods, *SQLearner* agents have close to zero variance across runs.

Binary-Cross-Entropy (BCE) loss for classification and the *mean-squared error (MSE)* loss for regression. We use the *Adam* (Kingma & Ba, 2014) optimizer (learning rate 0.003).

Second, the oracle network (Step 4, Section 2.4), that predicts an action for an input state. We encode the input state using 2 convolution layers with kernel-size 3 and *relu* activation. To predict the action, we use 3 fully-connected layers with *relu* activation and the BCE loss. We use *L2* regularization, and *Gradient Descent* with the *Adam* optimizer (learning rate 0.001).

4. Results

4.1. Learning optimal policies in Iterated Matrix Games using SQLoss

Iterated Prisoner’s Dilemma (IPD): We train different learners to play the IPD game. Figure 4 shows the results. For all learners, agents initially defect and move towards an NDR of -2.0 . This initial bias towards defection is expected, since, for agents trained with random game-play episodes, the benefits of exploitation outweigh the costs of mutual defection. For Selfish Learner (SL) agents, the bias intensifies, and the agents converge to mutually harmful selfish behavior (NDR= -2.0). Lola-PG agents learn to predict each other’s behavior and therefore realize that defection is more likely to lead to mutual harm than selfish benefit. They subsequently move towards cooperation, but occasionally defect (NDR= -1.2). In contrast, *SQLearner* agents quickly realize the costs of defection, indicated by the small initial dip in the NDR curves. They subsequently move towards almost perfect cooperation, with an NDR of -1.0 . Finally, it is important to note that *SQLearner* agents have close to zero variance, unlike other methods where the variance in NDR across runs is significant.

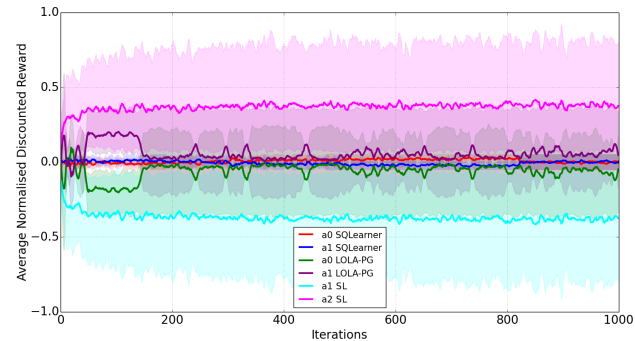


Figure 5: Average NDR values for different learners in the IMP game. *SQLearner* agents avoid exploitation by randomising between *H* and *T* to obtain a near-optimal NDR value (0) for this game. In contrast to other methods, *SQLearner* agents have close to zero variance across runs.

Iterated Matching Pennies (IMP): We train different learners to play the IMP game. The optimal policy for an agent to avoid exploitation is to play *H* or *T* perfectly randomly and obtain an NDR of 0. Figure 5 shows the results. *SQLearner* agents learn to play optimally and obtain an NDR close to 0. Interestingly, Selfish Learner and Lola-PG agents converge to an exploiter-exploited equilibrium where one agent consistently exploits the other agent. This asymmetric exploitation equilibrium is more pronounced for Selfish Learner agents than for Lola-PG agents. As before, we observe that *SQLearner* agents have close to zero variance across runs, unlike other methods where the variance in NDR across runs is significant.

Appendix B shows the results for the ISH game.

4.2. Evolving Cooperation in Games with visual input using GameDistill followed by SQLoss

4.2.1. THE COIN GAME: GAMEDISTILL

To evaluate the clustering step in *GameDistill*, we make two t-SNE (Maaten & Hinton, 2008) plots of the 100-dimensional feature vector extracted from the last layer of the *GameDistill* network. The first plot colors each point (state sequence) by the rewards obtained by both agents in the sequence. The second plot colors each point by the cluster label output by Agglomerative clustering. Figure 6 shows the results. *GameDistill* correctly learns two clusters, one for state sequences that represent cooperation and the other for state sequences that represent defection. We also experiment with different values for feature vector dimensions and obtain similar clustering results with sufficient training. Once we have the clusters, we train oracle networks using the state sequences in each cluster. To evaluate that the trained oracles represent a cooperation and a defection policy, we modify the Coin Game environment to contain only

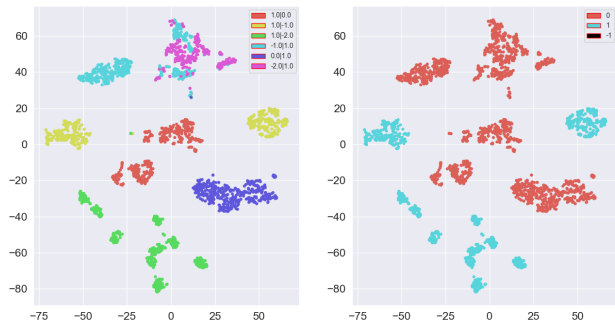


Figure 6: Representation of the clusters learned by *GameDistill*. Each point is a t-SNE projection of the 100-dimensional feature vector output by the *GameDistill* network for an input sequence of states. The figure on the left is colored based on actual rewards obtained by each agent (Red agent followed by Blue agent). The figure on the right is colored based on clusters learned by *GameDistill*. *GameDistill* correctly identifies two types of state sequences, one for cooperation (blue cluster) and the other for defection (red cluster).

the Red agent. We then play two variations of the game. In the first variation, the Red agent is forced to play the action suggested by the first oracle. In this variation, we find that the Red agent picks only 8.4% of Blue coins, indicating a high rate of cooperation. Therefore, the first oracle represents a cooperation policy. In the second variation, the Red agent is forced to play the action suggested by the second oracle. In this case, we find that the Red agent picks 99.4% of Blue coins, indicating a high rate of defection. Hence, the second oracle represents a defection policy. Therefore, the oracles learned by the Red agent using *GameDistill* represent cooperation and defection policies.

4.2.2. THE COIN GAME: SQLLOSS

Before playing the game, each *SQLearner* agent uses *GameDistill* to learn cooperation and defection oracles. During game-play, at each step, an agent follows either the action suggested by its cooperation oracle or the action suggested by its defection oracle. Further, each *SQLearner* agent has an additional *SQLoss* term. We compare approaches using the degree of cooperation between agents, measured by the probability of an agent to pick the coin of its color (Foerster et al., 2018). Figure 7 shows the results. The probability that an *SQLearner* agent will pick the coin of its color is close to 1. This high probability indicates that the other *SQLearner* agent is cooperating with this agent and only picking coins of its color. In contrast, the probability that a Lola-PG agent will pick its own coin is much smaller, indicating higher rates of defection. As expected, the probability of an agent picking its own coin is

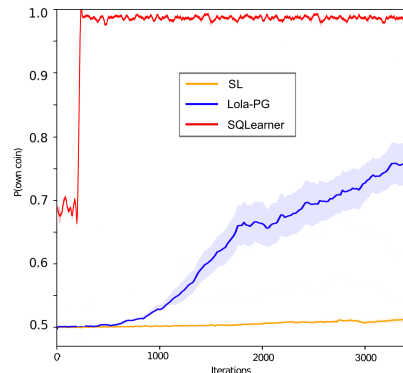


Figure 7: Probability of an agent picking a coin of its color for learners trained in the Coin Game. *SQLearner* agents cooperate (pick only their own coins) to achieve a near optimal strategy in the game. In contrast to Lola-PG, *SQLearner* agents have close to zero variance across runs.

the smallest for selfish learners (SL). The probability value of 0.5 indicates that a selfish learner is just as likely to pick the other agent’s coin as it is to pick its own coin.

4.3. SQLearner: Exploitability and Adaptability

Given that an agent does not have any prior information about the other agent, it is important that it evolves its strategy based on the strategy of its opponent. To evaluate an *SQLearner* agent’s ability to avoid exploitation by a selfish agent, we train one *SQLearner* agent against an agent that always defects in the Coin Game. We find that the *SQLearner* agent also learns to always defect. This persistent defection is important since given that the other agent is selfish, the *SQLearner* agent can do no better than also be selfish. To evaluate an *SQLearner* agent’s ability to exploit a cooperative agent, we train one *SQLearner* agent with an agent that always cooperates in the Coin Game. In this case, we find that the *SQLearner* agent learns to always defect. This persistent defection is important since given that the other agent is cooperative, the *SQLearner* agent obtains maximum reward by behaving selfishly. Hence, the *SQLearner* agent is both resistant to exploitation and able to exploit, depending on the strategy of the other agent.

5. Conclusion

We have described a status-quo loss (*SQLoss*) that encourages an agent to imagine the consequences of sticking to the status-quo. We demonstrated how agents trained with *SQLoss* evolve cooperative behavior in several social dilemmas without sharing rewards, gradients, or using a communication channel. To work with visual input games, we proposed *GameDistill*, an approach that automatically extracts a cooperative and a selfish policy from a social

dilemma game. We combined *GameDistill* and *SQLoss* to demonstrate how agents evolve desirable cooperative behavior in a social dilemma game with visual observations.

References

- Multiagent learning using a variable learning rate. *Artificial Intelligence*, 136(2):215–250, April 2002. ISSN 0004-3702.
- Abreu, D., Pearce, D., and Stacchetti, E. Toward a theory of discounted repeated games with imperfect monitoring. *Econometrica*, 58(5):1041–1063, 1990. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/2938299>.
- Axelrod, R. *Robert Axelrod's (1984) The Evolution of Cooperation*. Basic Books, 1984.
- Banerjee, D. and Sen, S. Reaching pareto-optimality in prisoner's dilemma using conditional joint action learning. *Autonomous Agents and Multi-Agent Systems*, 15(1), August 2007. ISSN 1387-2532.
- Code, A. Marl with sqloss. <https://github.com/user12423/MARL-with-SQLoss/>, 2019.
- Damer, S. and Gini, M. Achieving cooperation in a minimally constrained environment. volume 1, pp. 57–62, 01 2008.
- de Cote, E. M., Lazaric, A., and Restelli, M. Learning to cooperate in multi-agent social dilemmas. In *Proceedings of the Fifth International Joint Conference on Autonomous Agents and Multiagent Systems*, AAMAS '06, 2006.
- Dietz, T., Ostrom, E., and Stern, P. C. The struggle to govern the commons. *Science*, 302(5652):1907–1912, 2003. doi: 10.1126/science.1091015.
- Foerster, J., Chen, R. Y., Al-Shedivat, M., Whiteson, S., Abbeel, P., and Mordatch, I. Learning with opponent-learning awareness. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 122–130. International Foundation for Autonomous Agents and Multiagent Systems, 2018.
- Fudenberg, D. and Maskin, E. The folk theorem in repeated games with discounting or with incomplete information. *Econometrica*, 54(3):533–554, 1986. ISSN 00129682, 14680262.
- Fudenberg, D., Levine, D., and Maskin, E. The folk theorem with imperfect public information. *Econometrica*, 62(5):997–1039, 1994. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/2951505>.
- Green, E. J. and Porter, R. H. Noncooperative Collusion under Imperfect Price Information. *Econometrica*, 52(1): 87–100, 1984.
- Hardin, G. The tragedy of the commons. *Science*, 162(3859):1243–1248, 1968. doi: 10.1126/science.162.3859.1243.
- Hughes, E., Leibo, J. Z., Phillips, M., Tuyls, K., Dueñez Guzman, E., Castañeda, A. G., Dunning, I., Zhu, T., McKee, K., Koster, R., Roff, H., and Graepel, T. Inequity aversion improves cooperation in intertemporal social dilemmas. In *Proceedings of the 32Nd International Conference on Neural Information Processing Systems*, NIPS'18, 2018.
- Julien, P., Leibo, J., Zambaldi, V., Beattie, C., Tuyls, K., and Graepel, T. A multi-agent reinforcement learning model of common-pool resource appropriation. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, 12 2017.
- Kamada, Y. and Kominers, S. Information can wreck cooperation: A counterpoint to kandori (1992). *Economics Letters*, 107:112–114, 05 2010. doi: 10.1016/j.econlet.2009.12.040.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kleiman-Weiner, M., Ho, M. K., Austerweil, J. L., Littman, M. L., and Tenenbaum, J. B. Coordinate to cooperate or compete: abstract goals and joint intentions in social interaction. In *CogSci*, 2016.
- Lee, K. and Louis, K. *The Application of Decision Theory and Dynamic Programming to Adaptive Control Systems*. PhD thesis, 1967.
- Leibo, J. Z., Zambaldi, V., Lanctot, M., Marecki, J., and Graepel, T. Multi-agent reinforcement learning in sequential social dilemmas. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '17. International Foundation for Autonomous Agents and Multiagent Systems, 2017.
- Lerer, A. and Peysakhovich, A. Maintaining cooperation in complex social dilemmas using deep reinforcement learning, 2017.
- Luce, R. D. and Raiffa, H. *Games and decisions: Introduction and critical survey*. Courier Corporation, 1989.
- Maaten, L. v. d. and Hinton, G. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov): 2579–2605, 2008.

- Macy, D. and Flache, A. Learning dynamics in social dilemmas. *Proceedings of the National Academy of Sciences of the United States of America*, 99 Suppl 3:7229–36, 06 2002. doi: 10.1073/pnas.092080099.
- Nowak, M. and Sigmund, K. A strategy of win-stay, lose-shift that outperforms tit-for-tat in the prisoner’s dilemma game. *Nature*, 364:56–8, 08 1993. doi: 10.1038/364056a0.
- Nowak, M. A. and Sigmund, K. Tit for tat in heterogeneous populations. *Nature*, 355(6357):250–253, 1992.
- Nowak, M. A. and Sigmund, K. Evolution of indirect reciprocity by image scoring. *Nature*, 393(6685):573–577, 1998.
- Ohtsuki, H., Hauert, C., Lieberman, E., and Nowak, M. A. A simple rule for the evolution of cooperation on graphs and social networks. *Nature*, 441(7092):502–505, 2006. ISSN 1476-4687. doi: 10.1038/nature04605. URL <https://doi.org/10.1038/nature04605>.
- Ostrom, E. *Governing the commons-The evolution of institutions for collective actions*. Political economy of institutions and decisions, 1990.
- Ostrom, E., Burger, J., Field, C. B., Norgaard, R. B., and Policansky, D. Revisiting the commons: Local lessons, global challenges. *Science*, 284(5412):278–282, 1999. doi: 10.1126/science.284.5412.278.
- Peysakhovich, A. and Lerer, A. Consequentialist conditional cooperation in social dilemmas with imperfect information. In *International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018.
- Pineau, J. The Machine Learning Reproducibility Checklist. <https://www.cs.mcgill.ca/~jpineau/ReproducibilityChecklist.pdf>, 2019. URL <https://www.cs.mcgill.ca/~jpineau/ReproducibilityChecklist.pdf>.
- Sandholm, T. W. and Crites, R. H. Multiagent reinforcement learning in the iterated prisoner’s dilemma. *Bio Systems*, 37 1-2:147–66, 1996.
- Santos, F. and Pacheco, J. A new route to the evolution of cooperation. *Journal of evolutionary biology*, 19:726–33, 06 2006. doi: 10.1111/j.1420-9101.2005.01063.x.
- Sutton, R. S. and Barto, A. G. Reinforcement learning: An introduction. 2011.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pp. 1057–1063, 2000.
- Trivers, R. The evolution of reciprocal altruism. *Quarterly Review of Biology*, 46:35–57., 03 1971. doi: 10.1086/406755.
- Wang, J. X., Hughes, E., Fernando, C., Czarnecki, W. M., Duéñez Guzmán, E. A., and Leibo, J. Z. Evolving intrinsic motivations for altruistic behavior. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS ’19*, pp. 683–692. International Foundation for Autonomous Agents and Multiagent Systems, 2019. ISBN 978-1-4503-6309-9.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Wunder, M., Littman, M., and Babes, M. Classes of multiagent q-learning dynamics with ϵ -greedy exploration. In *Proceedings of the 27th International Conference on International Conference on Machine Learning, ICML’10*, 2010.
- Yu, C., Zhang, M., Ren, F., and Tan, G. Emotional multi-agent reinforcement learning in spatial social dilemmas. *IEEE Transactions on Neural Networks and Learning Systems*, 26(12):3083–3096, 2015.

Supplementary Material

A. Illustrations of Trained Oracle Networks for the Coin Game

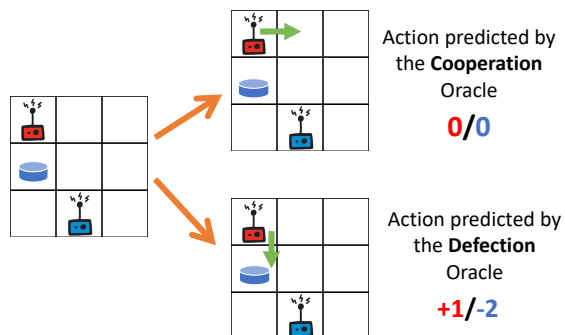


Figure 8: Illustrative predictions of the oracle networks learned by the Red agent using *GameDistill* in the Coin Game. The cooperation oracle suggests an action that avoids picking the coin of the other agent. The defection oracle suggests an action that involves picking the coin of the other agent.

Figure 8 shows the predictions of the oracle networks learned by the Red agent using *GameDistill* in the Coin Game. We see that the cooperation oracle suggests an action that avoids picking the coin of the other agent (the Blue coin). Analogously, the defection oracle suggests a selfish action that picks the coin of the other agent.

B. Results for the Iterated Stag Hunt using $SQLoss$

Figure 9 shows the results of training two *SQLearner* agents on the Iterated Stag Hunt game. *SQLearner* agents coordinate successfully to obtain a near-optimal NDR value (0) for this game.

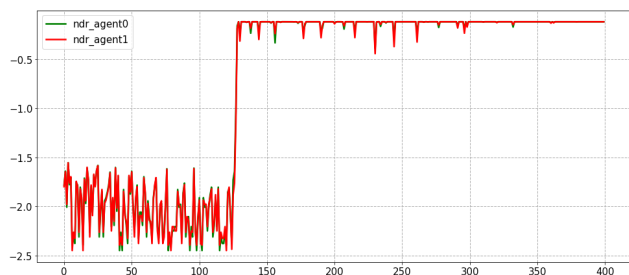


Figure 9: NDR values for *SQLearner* agents in the ISH game. *SQLearner* agents coordinate successfully to obtain a near optimal NDR value (0) for this game.

C. $SQLoss$: Effect of z on convergence to cooperation

We explore the effect of the hyper-parameter z (Section 2) on convergence to cooperation. To imagine the consequences of maintaining the status-quo, each agent samples κ_t from the Discrete Uniform distribution $\mathbb{U}\{1, z\}$. Therefore, a larger value of z implies a larger value of κ_t and longer imaginary episodes. We find that larger z (and hence κ) leads to faster cooperation between agents in the IPD and Coin Game. This effect plateaus at $z = 10$, which we select for our experiments.

D. Architecture Details

We performed all our experiments on an AWS instance with the following specifications.

- Model name: Intel(R) Xeon(R) Platinum 8275CL CPU @ 3.00GHz
- RAM: 189GB
- CPU(s): 96
- Architecture: x86_64
- Thread(s) per core: 2

E. Reproducibility Checklist

We follow the reproducibility checklist from (Pineau, 2019) and include further details here. For all the models and algorithms we have included details that we think would be useful for reproducing the results of this work.

- For all **models** and **algorithms** presented, check if you include:
 1. *A clear description of the mathematical setting, algorithm, and/or model:* **Yes**. The algorithm is described in detail in Section 2, with all the loss functions used for training being clearly defined. The details of the architecture, hyperparameters used and other algorithm details are given in Section 3. Environment details are explained in the sections that they are introduced in.
 2. *An analysis of the complexity (time, space, sample size) of any algorithm:* **No**. We do not include a formal complexity analysis of our algorithm. However, we do highlight the additional computational steps (in terms of losses and parameter updates) in Section 2 over standard multi-agent independently learning RL algorithms that would be needed in our approach.

3. *A link to a downloadable source code, with specification of all dependencies, including external libraries.:* **Yes.** We have made the source code available at [Code \(2019\)](#).
- For any **theoretical claim**, check if you include:
 1. *A statement of the result:* **NA.** Our paper is primarily empirical and we do not have any major theoretical claims. Hence this is Not Applicable.
 2. *A clear explanation of any assumptions:* **NA.**
 3. *A complete proof of the claim:* **NA.**
 - For all **figures** and **tables** that present empirical results, check if you include:
 1. *A complete description of the data collection process, including sample size:* **NA.** We did not collect any data for our work.
 2. *A link to a downloadable version of the dataset or simulation environment:* **Yes.** We have made the source code available at [Code \(2019\)](#).
 3. *An explanation of any data that were excluded, description of any pre-processing step:* **NA.** We did not perform any pre-processing step.
 4. *An explanation of how samples were allocated for training / validation / testing:* **Yes.** For *GameDistill* the details regarding data used for training is given in Section 2.4. The number of iterations used for learning (training) by *SQLearner* is given in Figures 4, 5 and 7. The details of the number of runs and the batch sizes used for various experiments are given in Section 3.
 5. *The range of hyper-parameters considered, method to select the best hyper-parameter configuration, and specification of all hyper-parameters used to generate results:* **Yes.** We did not do any hyperparameter tuning as part of this work. All the hyperparameters that we used are specified in Section 3.
 6. *The exact number of evaluation runs:* **Yes.** For all our environments, we repeat the experiment 20 times. For evaluation of performance, we use an average of 200 Monte Carlo estimates. We state this in Section 3. We do not fix any seeds. The details of the number of runs and the batch sizes used for various experiments are also given here.
 7. *A description of how experiments were run:* **Yes.** The README with instructions on how to run the experiments along with the source code is provided at [Code \(2019\)](#).
 8. *A clear definition of the specific measure or statistics used to report results:* **Yes.** We plot the mean and the one standard deviation region over the mean for all our numerical experiments. This is stated in Section 3.
 9. *Clearly defined error bars:* **Yes.** We plot the mean and the one standard deviation region over the mean for all our numerical experiments. This is stated in Section 3.
 10. *A description of results with central tendency (e.g. mean) & variation (e.g. stddev):* **Yes.** We plot the mean and the one standard deviation region over the mean for all our numerical experiments. This is stated in Section 3.
 11. *A description of the computing infrastructure used:* **Yes.** We have provided this detail in the Supplementary material in Section D.