
Inequity aversion improves cooperation in intertemporal social dilemmas

Edward Hughes*, Joel Z. Leibo*, Matthew Phillips, Karl Tuyls, Edgar Dueñez-Guzman,
Antonio García Castañeda, Iain Dunning, Tina Zhu, Kevin McKee, Raphael Koster,
Heather Roff, Thore Graepel

DeepMind, London, United Kingdom

{edwardhughes, jzl, karltuyls, duenez, antoniogc, idunning, tinazhu,
kevinrmckee, rkoster, hroff, thore}@google.com,
matthew.phillips.12@ucl.ac.uk

Abstract

Groups of humans are often able to find ways to cooperate with one another in complex, temporally extended social dilemmas. Models based on behavioral economics are only able to explain this phenomenon for unrealistic stateless matrix games. Recently, multi-agent reinforcement learning has been applied to generalize social dilemma problems to temporally and spatially extended Markov games. However, this has not yet generated an agent that learns to cooperate in social dilemmas as humans do. A key insight is that many, but not all, human individuals have inequity averse social preferences. This promotes a particular resolution of the matrix game social dilemma wherein inequity-averse individuals are personally pro-social and punish defectors. Here we extend this idea to Markov games and show that it promotes cooperation in several types of sequential social dilemma, via a profitable interaction with policy learnability. In particular, we find that inequity aversion improves temporal credit assignment for the important class of *intertemporal* social dilemmas. These results help explain how large-scale cooperation may emerge and persist.

1 Introduction

In intertemporal social dilemmas, there is a tradeoff between short-term individual incentives and long-term collective interest. Humans face such dilemmas when contributing to a collective food storage during the summer in preparation for a harsh winter, organizing annual maintenance of irrigation systems, or sustainably sharing a local fishery. Classical models of human behavior based on rational choice theory predict that cooperation in these situations is impossible [1, 2]. This poses a puzzle since humans evidently do find ways to cooperate in many everyday intertemporal social dilemmas, as documented by decades of fieldwork [3, 4] and laboratory experiments [5, 6]. Providing an empirically grounded explanation of how individual behavior gives rise to societal cooperation is seen as a core goal in several subfields of the social sciences and evolutionary biology [7, 8, 9].

[10, 11] proposed influential models based on behavioral game theory. However, these models have limited applicability since they only generate predictions when the problem can be cast as a matrix game (see e.g. [12, 13]). Here we consider a more realistic video-game setting, like those introduced in the behavioral research of [14, 15, 16]. In this environment, agents do not simply choose to cooperate or defect like they do in matrix games. Rather they must learn policies to implement their strategic decisions, and must do so while coping with the non-stationarity arising from other agents learning simultaneously. Several papers used multi-agent reinforcement learning [17, 18, 19] and

*Equal contribution.

planning [20, 21, 22, 23] to generate cooperation in this setting. However, this approach has not yet demonstrated robust cooperation in games with more than two players, which is often observed in human behavioral experiments. Moreover naïvely optimizing group reward is also ineffective, due to the lazy agent problem [24].[†]

It is difficult for both natural and artificial agents to find cooperative solutions to intertemporal social dilemmas for the following reasons:

1. Collective action – individuals must learn and coordinate policies at a group level to avoid falling into socially deficient equilibria.
2. Temporal credit assignment – rational defection in the short-term must become associated with long-term negative consequences.

Many different research traditions, including economics, evolutionary biology, sociology, psychology, and political philosophy have all converged on the idea that fairness norms are involved in resolving social dilemmas [25, 26, 27, 28, 29, 30, 31]. In one well-known model, agents are assumed to have inequity-averse preferences [10]. They balance their selfish desire for individual rewards against a need to keep deviations between their own rewards and the rewards of others as small as possible. Inequity-averse individuals are able to solve social dilemmas by resisting the temptation to pull ahead of others or—if punishment is possible—by punishing and discouraging free-riding. The inequity aversion model has been successfully applied to explain human behavior in a variety of laboratory economic games, such as the ultimatum game, the dictator game, the gift exchange game, market games, the trust game and public goods [32, 33].[‡]

In this research, we generalize the inequity aversion model to Markov games, and show that it resolves intertemporal social dilemmas. Crucial to our analysis will be the distinction between *disadvantageous* inequity aversion (negative reward received by individuals who underperform relative to others) and *advantageous* inequity aversion (negative reward received by individuals who overperform relative to others). Colloquially, these may be thought of as reductionist models of envy (disadvantageous inequity aversion) and guilt (advantageous inequity aversion) respectively [36]. We hypothesize that these directly address the two challenges set out above in the following way.

Inequity aversion mitigates the problem of collective action by changing the effective payoff structure experienced by agents through both a direct and an indirect mechanism. In the direct mechanism, defectors experience advantageous inequity aversion, diminishing the marginal benefit of defection over cooperation. The indirect mechanism arises when cooperating agents are disadvantageous-inequity averse. This motivates them to punish defectors by sanctioning them, reducing the payoff incentive for free-riding. Since agents must learn a defecting strategy via exploration, initially cooperative agents are deterred from switching strategies if the payoff bonus does not outweigh the cost of inefficiently executing the defecting strategy while learning.

Inequity aversion also ameliorates the temporal credit assignment problem. Learning the association between short-term actions and long-term consequences is a high-variance and error-prone process, both for animals [37] and reinforcement learning algorithms [38]. Inequity aversion short-circuits the need for such long-term temporal credit assignment by acting as an “early warning system” for intertemporal social dilemmas. As before, both a direct and an indirect mechanism are at work. With the direct mechanism, advantageous-inequity-averse defectors receive negative rewards in the short-term, since the benefits of defection are delivered on that timescale. The indirect mechanism operates because cooperators experience disadvantageous inequity aversion at precisely the time when other agents defect. This leads cooperators to punish defectors on a short-term timescale. Both systems have the effect of operant conditioning [39], incentivizing agents that cannot resolve long-term uncertainty to act in the lasting interest of the group.

2 Reinforcement learning in sequential social dilemmas

2.1 Partially observable Markov games

We consider multi-agent reinforcement learning in partially-observable general-sum Markov games [40, 41]. In each game state, agents take actions based on a partial observation of the state space and

[†]For more detail on the motivations for our research program, see the supplementary information.

[‡]For alternative theories of the other-regarding preferences that may underlie human cooperative behavior in economic games, see [34, 35].

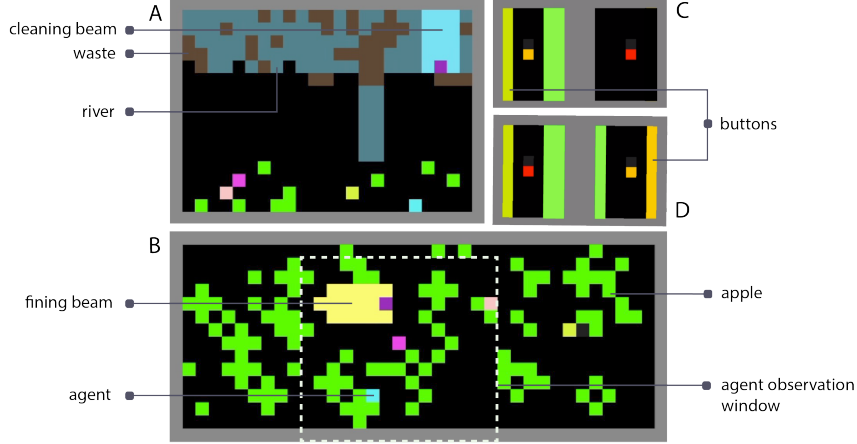


Figure 1: Screenshots from (A) the Cleanup game, (B) the Harvest game, (C) the Dictate apples game, and (D) the Take apples and Give apples games. The size of the agent-centered observation window is also shown in (B). The same size observation was used in all experiments.

receive an individual reward. Agents must learn through experience an appropriate behavior policy while interacting with one another. We formalize this as follows.

Consider an N -player partially observable Markov game \mathcal{M} defined on a finite set of states \mathcal{S} . The observation function $O : \mathcal{S} \times \{1, \dots, N\} \rightarrow \mathbb{R}^d$ specifies each player’s d -dimensional view on the state space. From each state, players may take actions from the set $\mathcal{A}^1, \dots, \mathcal{A}^N$ (one for each player). As a result of their joint action $a^1, \dots, a^N \in \mathcal{A}^1, \dots, \mathcal{A}^N$ the state changes following the stochastic transition function $\mathcal{T} : \mathcal{S} \times \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \Delta(\mathcal{S})$ (where $\Delta(\mathcal{S})$ denotes the set of discrete probability distributions over \mathcal{S}). Write $\mathcal{O}^i = \{o^i \mid s \in \mathcal{S}, o^i = O(s, i)\}$ to indicate the observation space of player i . Each player receives an individual extrinsic reward defined as $r^i : \mathcal{S} \times \mathcal{A}^1 \times \dots \times \mathcal{A}^N \rightarrow \mathbb{R}$ for player i .[§]

Each agent learns, independently through its own experience of the environment, a behavior policy $\pi^i : \mathcal{O}^i \rightarrow \Delta(\mathcal{A}^i)$ (written $\pi(a^i | o^i)$) based on its own observation $o^i = O(s, i)$ and extrinsic reward $r^i(s, a^1, \dots, a^N)$. For the sake of simplicity we will write $\vec{a} = (a^1, \dots, a^N)$, $\vec{o} = (o^1, \dots, o^N)$ and $\vec{\pi}(\cdot | \vec{o}) = (\pi^1(\cdot | o^1), \dots, \pi^N(\cdot | o^N))$. Each agent’s goal is to maximize a long term γ -discounted payoff defined as follows:

$$V_{\vec{\pi}}^i(s_0) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r^i(s_t, \vec{a}_t) \mid \vec{a}_t \sim \vec{\pi}_t, s_{t+1} \sim \mathcal{T}(s_t, \vec{a}_t) \right]. \quad (1)$$

2.2 Learning agents

We deploy asynchronous advantage actor-critic (A3C) as the learning algorithm for our agents [42]. A3C maintains both value (critic) and policy (actor) estimates using a deep neural network. The policy is updated according to the policy gradient method, using a value estimate as a baseline to reduce variance. Gradients are generated asynchronously by 24 independent copies of each agent, playing simultaneously in distinct instantiations of the environment. Explicitly, the gradients are $\nabla_{\theta} \log \pi(a_t | s_t; \theta) A(s_t, a_t; \theta, \theta_v)$, where $A(s_t, a_t; \theta, \theta_v)$ is the advantage function, estimated via k -step backups, $\sum_{i=0}^{k-1} \gamma^i u_{t+i} + \gamma^k V(s_{t+k}; \theta_v) - V(s_t; \theta_v)$ where u_{t+i} is the *subjective* reward. In section 3.1 we decompose this into an *extrinsic* reward from the environment and an *intrinsic* reward that defines the agent’s inequity-aversion.

2.3 Intertemporal social dilemmas

An intertemporal social dilemma is a temporally extended multi-agent game in which individual short-term optimal strategies lead to poor long-term outcomes for the group. To define this term

[§]In our games, $N = 5$, $d = 15 \times 15 \times 3$ and $|\mathcal{A}^i|$ ranges from 8 to 10, with actions comprising movement, rotation and firing.

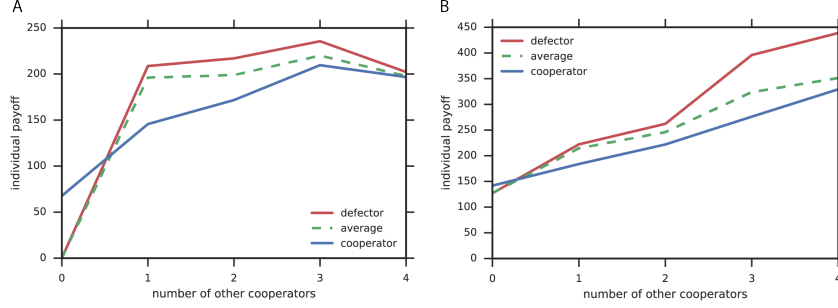


Figure 2: The public goods game (Cleanup) and the commons game (Harvest) are social dilemmas. (A) shows the Schelling diagram for Cleanup. (B) shows the Schelling diagram for Harvest. The dotted line shows the overall average return were the individual to choose defection.

precisely, we employ a formalization of empirical game theoretic analysis [43, 44]. Our definition is consistent with that of [17]. However, since that work was limited to the 2-player case, it relied on the empirical payoff matrix to represent the relative values of cooperation and defection. This quantity is unwieldy for $N > 2$ since it becomes a tensor. Therefore we base our definition on a different representation of the N -player game. Explicitly, a *Schelling diagram* [45, 18] depicts the relative payoffs for a single cooperator or defector given a fixed number of other cooperators. Thus Schelling diagrams are a natural and convenient generalization of payoff matrices to multi-agent settings. Game-theoretic properties like Nash equilibria are readily visible in Schelling diagrams; see [45] for additional details and intuition.

An N -player *sequential social dilemma* is a tuple $(\mathcal{M}, \Pi = \Pi_c \sqcup \Pi_d)$ of a Markov game and two disjoint sets of policies, said to implement cooperation and defection respectively, satisfying the following properties. Consider the strategy profile $(\pi_c^1, \dots, \pi_c^\ell, \pi_d^1, \dots, \pi_d^m) \in \Pi_c^\ell \times \Pi_d^m$ with $\ell + m = N$. We shall denote the average payoff for the cooperating policies by $R_c(\ell)$ and for the defecting policies by $R_d(\ell)$. A *Schelling diagram* plots the curves $R_c(\ell + 1)$ and $R_d(\ell)$. Intuitively, the diagram displays the two possible payoffs to the N^{th} player given that ℓ of the remaining players elect to cooperate and the rest defect. We say that (\mathcal{M}, Π) is a sequential social dilemma iff the following hold:

1. Mutual cooperation is preferred over mutual defection: $R_c(N) > R_d(0)$.
2. Mutual cooperation is preferred to being exploited by defectors: $R_c(N) > R_c(0)$.
3. Either the *fear* property, the *greed* property, or both:
 - Fear: mutual defection is preferred to being exploited. $R_d(i) > R_c(i)$ for sufficiently small i .
 - Greed: exploiting a cooperator is preferred to mutual cooperation. $R_d(i) > R_c(i)$ for sufficiently large i .

We show that the matrix games Stag Hunt, Chicken and Prisoner’s Dilemma satisfy these properties in Supplementary Fig. 1.

A sequential social dilemma is *intertemporal* if the choice to defect is optimal in the short-term. More precisely, consider an individual i and an arbitrary set of policies for the rest of the group. Given a starting state, for all k sufficiently small, the policy $\pi_k^i \in \Pi$ with maximum return in the next k steps is a defecting policy. There is thus a tension between short-term personal gain and long-term group utility.

2.4 Examples

[46] divides all multi-person social dilemmas into two broad categories:

1. *Public goods dilemmas*, in which an individual must pay a personal cost in order to provide a resource that is shared by all.
2. *Commons dilemmas*, in which an individual is tempted by a personal benefit, depleting a resource that is shared by all.

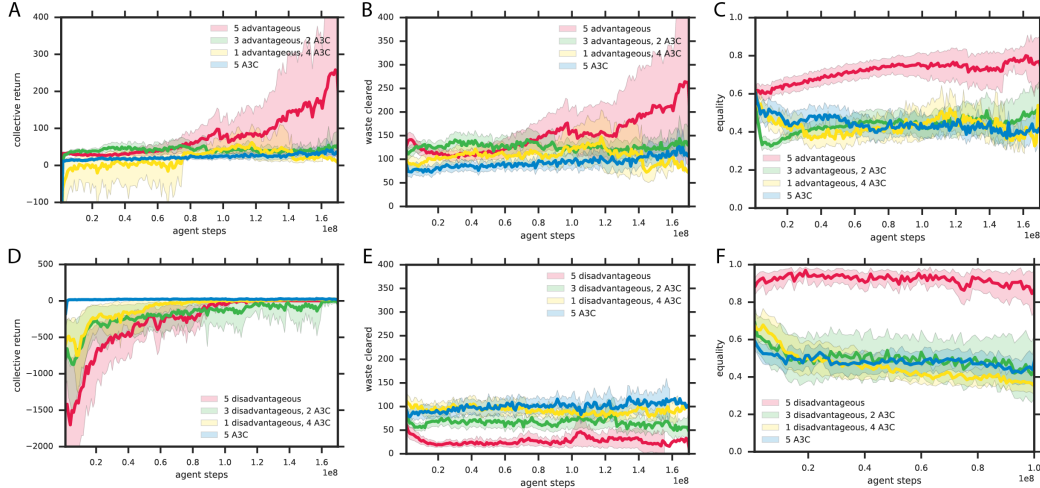


Figure 3: Advantageous inequity aversion facilitates cooperation in the Cleanup game. (A) compares the collective return achieved by A3C and advantageous inequity averse agents, (B) shows contributions to the public good, and (C) shows equality over the course of training. (D-F) demonstrate that disadvantaged inequity aversion does not promote greater cooperation in the Cleanup game.

We consider two dilemmas in this paper, one of the public goods type and one of the commons type. Each was implemented as a partially observable Markov game on a 2D grid. Both are also intertemporal social dilemmas because individually selfish actions produce immediate benefits while their impacts on the collective develop over a longer time horizon. The availability of costly punishment is of critical importance in human sequential social dilemmas [47, 48] and is therefore an action in the environments presented here.[¶]

In the *Cleanup* game, the aim is to collect apples from a field. Each apple provides a reward of 1. The spawning of apples is controlled by a geographically separate aquifer that supplies water and nutrients. Over time, this aquifer fills up with waste, lowering the respawn rate of apples linearly. For sufficiently high waste levels, no apples can spawn. At the start of each episode, the environment resets with waste just beyond this saturation point. To cause apples to spawn, agents must clean some of the waste.

Here we have a dilemma. Provided that some agents contribute to the public good by cleaning up the aquifer, it is individually more rewarding to stay in the apple field. However, if all players defect, then no-one gets any reward. A successful group must balance the temptation to free-ride with the provision of the public good. Cooperative agents must make a positive commitment to group-level well-being to solve the task.

The goal of the *Harvest* game is to collect apples. Each apple provides a reward of 1. The apple regrowth rate varies across the map, dependent on the spatial configuration of uncollected apples: the more nearby apples, the higher the local regrowth rate. If all apples in a local area are harvested then none ever grow back. After 1000 steps the episode ends, at which point the game resets to an initial state.

The dilemma is as follows. The short-term interests of each individual leads toward harvesting as rapidly as possible. However, the long-term interests of the group as a whole are advanced if individuals refrain from doing so, especially when many agents are in the same local region. Such situations are precarious because the more harvesting agents there are, the greater the chance of permanently depleting the local resources. Cooperators must abstain from a personal benefit for the good of the group.^{||}

[¶]In both games, players can fine each other using a punishment beam. This contrasts with [18], in which a timeout beam was used.

^{||}Precise details of the ecological dynamics may be found in the supplementary information.

2.5 Validating the environments

We would like to demonstrate that these environments are social dilemmas by plotting Schelling diagrams. In complex, spatially and temporally extended Markov games, it is not feasible to analytically determine cooperating and defecting policies. Instead, we must study the environment empirically. One method employs reinforcement learning to train such policies. We enforce cooperation or defection by making appropriate modifications to the environment, as follows.

In Harvest, we enforce cooperation by modifying the environment to prevent some agents from gathering apples in low-density areas. In Cleanup, we enforce free-riding by removing the ability of some agents to clean up waste. We also add a small group reward signal to encourage the remaining agents to cooperate. The resulting empirical Schelling diagrams in Figure 2 prove that our environments are indeed social dilemmas.

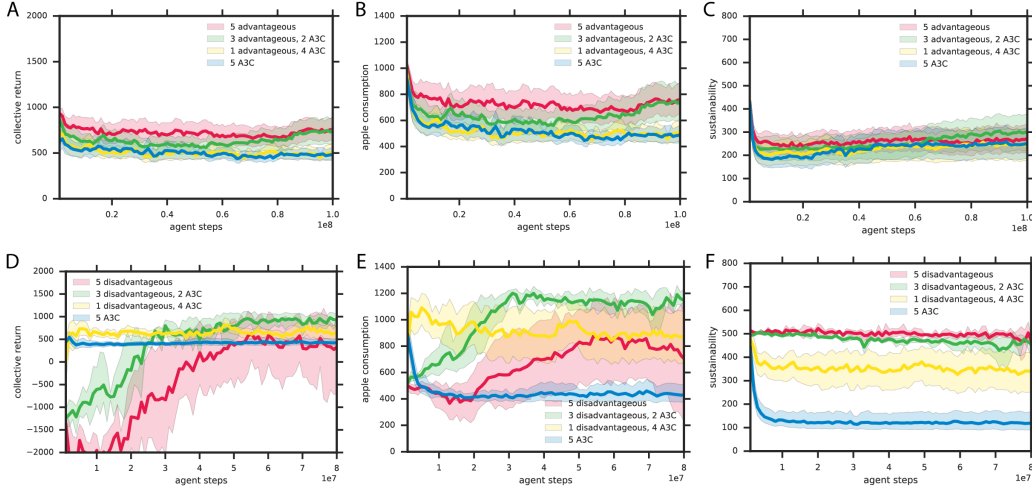


Figure 4: Inequity aversion promotes cooperation in the Harvest game. When all 5 agents have advantageous inequity aversion, there is a small improvement over A3C in the three social outcome metrics: (A) collective return, (B) apple consumption, and (C) sustainability. Disadvantageous inequity aversion provides a much larger improvement over A3C, and works even when only 1 out of 5 agents are inequity averse. (D) shows collective return, (E) apple consumption, and (F) sustainability.

3 The model

We first introduce the inequity aversion model of [10]. It is directly applicable only to stateless games. We then extend their model to sequential or multi-state problems, making use of deep reinforcement learning.

3.1 Inequity aversion

The [10] utility function is as follows. Let r_1, \dots, r_N be the extrinsic payoffs achieved by each of N players. Each agent receives a utility

$$U_i(r_i, \dots, r_N) = r_i - \frac{\alpha_i}{N-1} \sum_{j \neq i} \max(r_j - r_i, 0) - \frac{\beta_i}{N-1} \sum_{j \neq i} \max(r_i - r_j, 0), \quad (2)$$

where the additional terms may be interpreted as intrinsic payoffs, in the language of [49].

The parameter α_i controls an agent's aversion to *disadvantageous* inequity. A larger value for α_i implies a larger utility loss when other agents achieve rewards greater than one's own. Likewise, the parameter β_i controls an agent's aversion to *advantageous* inequity, utility lost when performing better than others. [10] argue that $\alpha > \beta$. That is, most people are loss averse in social comparisons. There is some empirical support for this prediction [50], though the evidence is mixed [51, 52]. In a sweep over values for α and β , we found our strongest results for $\alpha = 5$ and $\beta = 0.05$.

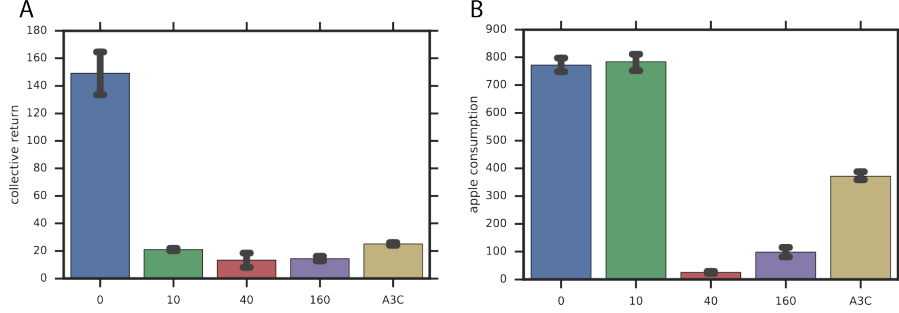


Figure 5: Inequity aversion promotes cooperation by improving temporal credit assignment. (A) shows collective return for delayed advantageous inequity aversion in the Cleanup game. (B) shows apple consumption for delayed disadvantageous inequity aversion in the Harvest game.

3.2 Inequity aversion in sequential dilemmas

Experimental work in behavioral economics suggests that some proportion of natural human populations are inequity averse [8]. However, as a computational model, inequity aversion has only been expounded for the matrix game setting. Equation (2) can be directly applied only to stateless games [53, 54]. In this section we extend this model of inequity aversion to the temporally extended Markov game case.

The main problem in re-defining the social preference of equation (2) for Markov games is that the rewards of different players may occur on different timesteps. Thus the key step in extending (2) to this case is to introduce per-player temporal smoothing of the reward traces.

Let $r_i(s, a)$ denote the reward obtained by the i -th player when it takes action a from state s . For convenience, we also sometimes write it with a time index: $r_i^t := r_i(s^t, a^t)$. We define the subjective reward $u_i(s, a)$ received by the i -th player when it takes action a from state s to be

$$u_i(s_i^t, a_i^t) = r_i(s_i^t, a_i^t) - \frac{\alpha_i}{N-1} \sum_{j \neq i} \max(e_j^t(s_j^t, a_j^t) - e_i^t(s_i^t, a_i^t), 0) - \frac{\beta_i}{N-1} \sum_{j \neq i} \max(e_j^t(s_j^t, a_j^t) - e_j^t(s_j^t, a_j^t), 0), \quad (3)$$

where the temporal smoothed rewards e_j^t for the agents $j = 1, \dots, N$ are updated at each timestep t according to

$$e_j^t(s_j^t, a_j^t) = \gamma \lambda e_j^{t-1}(s_j^{t-1}, a_j^{t-1}) + r_j^t(s_j^t, a_j^t), \quad (4)$$

where γ is the discount factor and λ is a hyperparameter. This is analogous to the mathematical formalism used for eligibility traces [55]. Furthermore, we allow agents to observe the smoothed reward of every player on each timestep.

4 Results

We show that advantageous inequity aversion is able to resolve certain intertemporal social dilemmas without resorting to punishment by providing a temporally correct intrinsic reward. For this mechanism to be effective, the population must have sufficiently many advantageous-inequity-averse individuals. By contrast disadvantageous-inequity-averse agents can drive mutual cooperation even in small numbers. They achieve this by punishing defectors at a time concomitant with their offences. In addition, we find that advantageous inequity aversion is particularly effective for resolving public goods dilemmas, whereas disadvantageous inequity aversion is more powerful for addressing commons dilemmas. Our baseline A3C agent fails to find socially beneficial outcomes in either category of game. We define the metrics used to quantify our results in the supplementary information.

4.1 Advantageous inequity aversion promotes cooperation

Advantageous-inequity-averse agents are better than A3C at maintaining cooperation in both public goods and commons games. This effect is particularly pronounced in the Cleanup game (Figure 3).

Here groups of 5 advantageous-inequity-averse agents find solutions in which 2 consistently clean large amounts of waste, producing a large collective return.** We clarify the effect of advantageous inequity aversion on the intertemporal nature of the problem by delaying the delivery of the intrinsic reward signal. Figure 5 suggests that improving temporal credit assignment is an important function of inequity aversion since delaying the time at which the intrinsic reward signal is delivered removes its beneficial effect.

4.2 Disadvantageous inequity aversion promotes cooperation

Disadvantageous-inequity-averse agents are better than A3C at maintaining cooperation via punishment in commons games (Figure 4). In particular, a single disadvantageous-averse agent can fine defectors, generating a sustainable outcome.†† In Figure 5, we see that the disadvantageous-inequity-aversion signal must be temporally aligned with over-consumption for effective policing to arise. Hence, it is plausible that inequity aversion bridges the temporal gap between short-term incentives and long-term outcomes. Disadvantageous inequity aversion has no such positive impact in the Cleanup game, for reasons that we discuss in section 5.

5 Discussion

In the Cleanup game, advantageous inequity aversion is an unambiguous feedback signal: it encourages agents to contribute to the public good. In the direct pathway, trial and error will quickly discover that the fastest way to diminish the negative rewards arising from advantageous inequity aversion is to clean up waste, since doing so creates more apples for others to consume. However the indirect mechanism of disadvantageous inequity aversion and punishment lacks this property; while punishment may help exploration of new policies, it does not directly increase the attractiveness of waste cleaning.

The Harvest game requires passive abstention rather than active provision. In this setting, advantageous inequity aversion provides a noisy signal for sustainable behaviour. This is because it is sensitive to the precise apple configuration in the environment, which changes rapidly over time. Hence advantageous inequity aversion does not greatly aid the exploration of policy space. Punishment, on the other hand, operates as a valuable shaping reward for learning, dis-incentivizing overconsumption at precisely the correct time and place.

In the Harvest game, disadvantageous inequity aversion generates cooperation in a grossly inefficient manner: huge amounts of collective resource are lost to fines (compare Figures 4D and 4E). This parallels human behavior in laboratory matrix games, e.g. [56, 57]. In the Cleanup game, advantageous-inequity averse agents resolve the social dilemma without such losses, but must comprise a large proportion of the population to be successful. This mirrors the cultural modulation of advantageous inequity aversion in humans [58]. Evolution is hypothesized to have favored fairness as a mechanism for continued human cooperation [59]. It remains to be seen whether emergent inequity-aversion can be obtained by evolving reinforcement learning agents.

We conclude by putting our approach in the context of prior work. Since our mechanism does not require explicitly training cooperating and defecting agents or modelling their behaviour, it scales more easily to complex environments and large populations of agents. However, our method has several limitations. Firstly, our guilty agents are quite exploitable, as evidenced by the necessity of a homogeneous guilty population to achieve cooperation. Secondly, our agents use outcomes rather than predictions to inform their policies. This is known to be a problem in environments with high stochasticity [22]. Finally, the heterogeneity of the population is an additional hyperparameter in our model. Clearly, one must set this appropriately, particularly in games with asymmetric outcomes. It is likely that a hybrid approach will be required to solve these challenging issues at scale.

**For a video of this behavior, visit <https://youtu.be/N8BUzzFx7uQ>.

††For a video of this behavior, visit <https://youtu.be/tz3ZpTTmxTk>.

A Supplementary information

A.1 Motivating research on emergent cooperation

The aims of this new research program are twofold. First, we seek to better understand the individual level inductive biases that promote emergent cooperation at the group level in humans. Second, we want to develop agents that exhibit these inductive biases, in the hope that they might navigate complex multi-agent tasks in a human-like way. Much as the fields of neuroscience and reinforcement learning have enjoyed a symbiotic relationship over the past fifty years, so also can behavioral economics and multi-agent reinforcement learning.

Consider, for comparison, maximizing joint utility. Firstly, this assumes away the problem of emergent altruism on the individual level, which is exactly our object of study. Therefore, it is not a relevant baseline for our research. Moreover, it is known to suffer from a serious spurious reward problem (Sunehag et al. 2017), which gets worse as the number of agents increases. Furthermore, in realistic environments, one may not have access to the collective reward function, for privacy reasons for example. Finally, groups of agents trained with a group reward are by definition overfitting to the outcomes of their co-players. Thus maximizing joint utility does not easily generalize to complicated multi-agent problems with large numbers of agents and subtasks that mix cooperation and competition.

Individual-level inductive biases sidestep these issues, while allowing us to learn from the extensive human behavioral literature. In this paper, we have taken an extremely well-studied model in the game-theoretic setting (Fehr and Schmidt 1999) and recast it as an intrinsic reward for reinforcement learning. We can thus evaluate the strengths and weaknesses of inequity aversion from a completely new perspective. We note its success in solving social dilemmas, but find that the success is task-conditional, and that the policies are sometimes quite exploitable. This suggests various fascinating extensions, such as a population-based study with evolved intrinsic rewards (Wang et al. to appear).

A.2 Illustrative Schelling diagrams for 2-player matrix games and SSDs

Figure 1 shows Schelling diagrams and the associated payoff matrices for the canonical matrix games Chicken, Stag Hunt and Prisoner’s Dilemma. We may read off the pure strategy Nash equilibria by considering the social pressure generated by the dominant strategy. Where this is defection, then there is a negative pressure on the number of cooperators; where this is cooperation, there is a positive pressure. Hence the pure strategy Nash equilibria in Chicken are (c, d) and (d, c) , in Stag Hunt (c, c) and (d, d) and in Prisoner’s Dilemma (d, d) . Moreover, the different motivations for defection are immediately apparent. In Chicken, greed promotes defection: $R_d(1) > R_c(1)$. In Stag Hunt, the problem is fear: $R_d(0) > R_c(0)$. Prisoner’s Dilemma suffers from both temptations to defect.

A.3 Parameters for Cleanup and Harvest games

In both Cleanup and Harvest, all agents are equipped with a fining beam which administers -1 reward to the user and -50 reward to the individual that is being fined. There is no penalty to the user for unsuccessful fining. In Cleanup each agent is additionally equipped with a cleaning beam, which allows them to remove waste from the aquifer. In both games, eating apples provides a reward of 1. There are no other extrinsic rewards.

In Cleanup, waste is produced uniformly in the river with probability 0.5 on each timestep, until the river is saturated with waste, which happens when the waste covers 40% of the river. For a given saturation x of the river, apples spawn in the field with probability $0.125x$. Initially the river is saturated with waste, so some contribution to the public good is required for any agent to receive a reward.

In Harvest, apples spawn relative to the current number of other apples within an ℓ^1 radius of 2. The spawn probabilities are 0, 0.005, 0.02, 0.05 for 0, 1, 2 and ≥ 3 apples inside the radius respectively. The initial distribution of apples creates a number of more or less precariously linked regions. Sustainable policies must preferentially harvest denser regions, and avoid removing the important apples that link patches.

A.4 Social outcome metrics

Unlike in single-agent reinforcement learning where the value function is the canonical metric of agent performance, in multi-agent systems with mixed incentives, there is no scalar metric that can



Figure 6: These Schelling diagrams demonstrate that classic matrix games are social dilemmas by our definition.

adequately track the state of the system (see e.g. [60, 18]). Thus we use several different social outcome metrics in order to summarize group behavior and facilitate its analysis.

Consider N independent agents. Let $\{r_t^i \mid t = 1, \dots, T\}$ be the sequence of rewards obtained by the i -th agent over an episode of duration T . Likewise, let $\{o_t^i \mid t = 1, \dots, T\}$ be the i -th agent's observation sequence. Its return is given by $R^i = \sum_{t=1}^T r_t^i$.

The *Utilitarian metric* (U), also known as *collective return*, measures the sum total of all rewards obtained by all agents. It is defined as the average over players of sum of rewards R^i . The *Equality metric* (E) is defined using the Gini coefficient [61]. The *Sustainability metric* (S) is defined as the average time at which the rewards are collected. For the Cleanup game, we also consider a measure of total contribution to the public good (P), defined as the number of waste cells cleaned.

$$U = \mathbb{E} \left[\frac{\sum_{i=1}^N R^i}{T} \right], \quad (5)$$

$$E = 1 - \frac{\sum_{i=1}^N \sum_{j=1}^N |R^i - R^j|}{2N \sum_{i=1}^N R^i}, \quad (6)$$

$$S = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N t^i \right] \quad \text{where } t^i = \mathbb{E}[t \mid r_t^i > 0], \quad (7)$$

$$P = \sum_i^N p_i. \quad (8)$$

where p_i is the number of waste cells cleaned by player i .

A.5 Dictate apples, Give apples and Take apples games

In each game, two players are isolated from one another in separate “rooms”. They can interact only by pressing buttons. In the *Dictate* apples game, initially all apples are in the left room. At any time, the left agent can press a button that transports all the apples it has not yet consumed to the right room. In the *Take* apples game, both players begin with apples in their room, but there are twice as many in the left room as the right room. The right agent has the option at any time of pressing a button that removes all the apples from the other player's room that have not yet been collected. In the *Give* apples game, both players begin with apples, and the left player again has twice as many as the right player. The left player can press a button to add more apples on the right side. Unlike in the *Dictate* apples game, this has no effect on the left agent's own apple supply. Each episode terminates when all apples are collected.

A.6 Inequity aversion models “irrational” behavior

The inequity aversion model of [10] is supported by experimental evidence from behavioral game theory. In particular, human behavior in the Dictator game is consistent with the prediction that some people have inequity-averse social preferences. A subject in a typical Dictator game experiment must decide how much of an initial endowment (if any) to give to another subject in a one-shot anonymous manner. In contrast to the prediction of rational choice theory that subjects would offer 0—but in accord with the prediction of [10]’s inequity aversion model—most subjects offer between 10% and 50% [62].

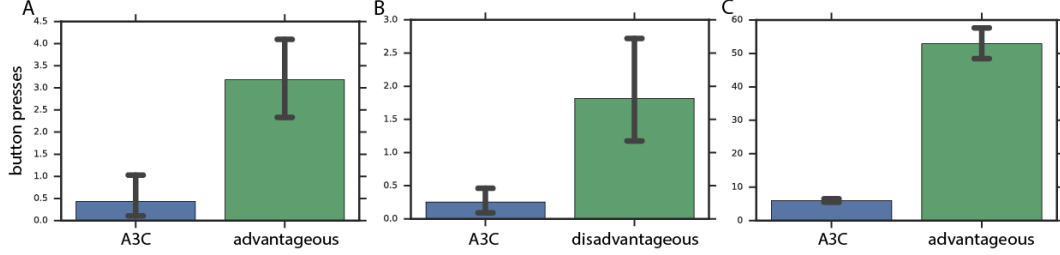


Figure 7: Behavioral economics laboratory paradigms can be simulated by gridworld Markov games. Agent behavior is shown in (A) for the Dictate apples game, in (B) for the Take apples game, and in (C) for the Give apples game.

To test whether our temporally extended inequity-aversion model makes predictions consistent with these findings, we introduce 3 simple 2-player gridworld games (see Figure 1). These capture the essential features of Dictator game laboratory experiments. As in all our experiments, positive agent external rewards can only be obtained by collecting apples. In addition an agent can press buttons which *Dictate* apples (give from its own store), *Give* apples from an external store or *Take* apples from the other agent. A full description is provided in the supplementary information.

A selfish rational agent would never press the button in any of these games. This prediction was borne out by our A3C agent baseline (Figure 7). On the other hand, advantageous-inequity-averse agents pressed their buttons significantly more often in the Give apples and Dictate apples games. They pressed the button even in the Dictate apples game when doing so could only reduce their own (extrinsic) payoff. Disadvantageous-inequity-averse agents pressed their button in the Take apples game to reduce the rewards obtained by the player with the larger initial endowment despite there being no extrinsic benefit to doing this.

A.7 Theoretical arguments for the success of inequity aversion

We provide theoretical arguments for inequity aversion as an improvement to temporal credit assignment, extending the work of (Fehr and Schmidt 1999) beyond simple market games. In an intertemporal social dilemma, defection dominates cooperation in the short term. To leading order, the short-term Schelling diagram for an intertemporal social dilemma looks like Figure 8A, since by definition defection must dominate cooperation. Here and in the sequel we work in the limit of large number of players N . Mathematically, we denote defector payoff by D , cooperator payoff by C and average payoff across the population by \bar{R} , writing:

$$C = c, \quad D = d, \quad \bar{R} = -\frac{d-c}{N}x + d, \quad \text{with } d > c. \quad (9)$$

First consider the effect of advantageous inequity aversion (AIA) on the short-term payoffs. Clearly the cooperator line is unchanged, since it is dominated. Hence the cooperator and defector lines become:

$$\tilde{C} = c, \quad \tilde{D} = D - \alpha(D - \bar{R}) = d - \alpha\frac{(d-c)}{N}x, \quad \text{with } \alpha > 0. \quad (10)$$

The transformed short-term payoffs are shown in Figure 8B. Since the C curve dominates \tilde{D} in some region, cooperative behavior can be self-sustaining in the short-term. Thus AIA improves temporal credit assignment. AIA can resolve the social dilemma when the earliest learned behavior generates *multiple* cooperators. This is the case for the Cleanup game but not the Harvest game, explaining the results.

The primary effect of disadvantageous inequity aversion (DIA) is to lower the payoff to a cooperator. However, it also motivates the cooperator to use the fining tool to reduce \bar{R} . There are several simple reasons why defectors might end up being especially targeted. Firstly, the behavior that avoids the policing agent may be cooperative (as in the Harvest game). Secondly, policing agents are motivated to avoid tagging other policers, because of the danger of retaliation.

Assuming that defectors are especially targeted, the cooperator and defector lines become:

$$\tilde{C} = C - \beta_C(\bar{R} - C) = c + \beta_C \left(\frac{d-c}{N}x - (d-c) \right), \quad (11)$$

$$\tilde{D} = D - \beta_D(\bar{R} - C) = d + \beta_D \left(\frac{d-c}{N}x - (d-c) \right), \quad (12)$$

with $\beta_D > \beta_C > 0$. The transformed short-term payoffs are shown in Figure 8C. Here the Nash equilibrium has moved to a positive number of cooperators. Hence DIA has improved temporal credit assignment. Of course, this argument requires the policing effect to emerge in the first place. This is possible when the earliest learned behavior is defection (Harvest), but not when it is cooperation (Cleanup), explaining the results.

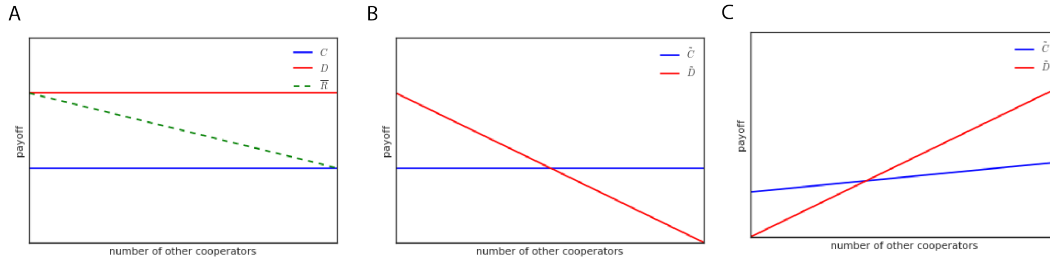


Figure 8: Inequity aversion alters the effective payoffs from cooperation and defection in the short-term, in such a way that cooperative behavior is rationally learnable. Hence, helps to solve the intertemporal social dilemma.

References

- [1] M. Olson, *The logic of collective action*. Harvard University Press, 1965.
- [2] G. Hardin, “The tragedy of the commons,” *Science*, vol. 162, no. 3859, pp. 1243–1248, 1968.
- [3] E. Ostrom, *Governing the Commons: The Evolution of Institutions for Collective Action*. Cambridge University Press, 1990.
- [4] T. Dietz, E. Ostrom, and P. C. Stern, “The struggle to govern the commons,” *science*, vol. 302, no. 5652, pp. 1907–1912, 2003.
- [5] E. Ostrom, J. Walker, and R. Gardner, “Covenants with and without a sword: Self-governance is possible,” *American political science Review*, vol. 86, no. 02, pp. 404–417, 1992.
- [6] E. Fehr and S. Gächter, “Altruistic punishment in humans,” *Nature*, vol. 415, no. 6868, p. 137, 2002.
- [7] E. Ostrom, “A behavioral approach to the rational choice theory of collective action: Presidential address, american political science association, 1997,” *American political science review*, vol. 92, no. 1, pp. 1–22, 1998.
- [8] E. Fehr and H. Gintis, “Human motivation and social cooperation: Experimental and analytical foundations,” *Annu. Rev. Sociol.*, vol. 33, pp. 43–64, 2007.
- [9] D. G. Rand and M. A. Nowak, “Human cooperation,” *Trends in cognitive sciences*, vol. 17, no. 8, pp. 413–425, 2013.
- [10] E. Fehr and K. M. Schmidt, “A theory of fairness, competition, and cooperation,” *The quarterly journal of economics*, vol. 114, no. 3, pp. 817–868, 1999.
- [11] A. Falk and U. Fischbacher, “A theory of reciprocity,” *Games and economic behavior*, vol. 54, no. 2, pp. 293–315, 2006.
- [12] T. W. Sandholm and R. H. Crites, “Multiagent reinforcement learning in the iterated prisoner’s dilemma,” *Biosystems*, vol. 37, no. 1-2, pp. 147–166, 1996.
- [13] E. Munoz de Cote, A. Lazaric, and M. Restelli, “Learning to cooperate in multi-agent social dilemmas,” in *5th International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS 2006)*, pp. 783–790, 2006.
- [14] M. A. Janssen, R. Holahan, A. Lee, and E. Ostrom, “Lab experiments for the study of social-ecological systems,” *Science*, vol. 328, no. 5978, pp. 613–617, 2010.
- [15] M. Janssen, “Introducing ecological dynamics into common-pool resource experiments,” *Ecology and Society*, vol. 15, no. 2, 2010.
- [16] M. Janssen, “The role of information in governing the commons: experimental results,” *Ecology and Society*, vol. 18, no. 4, 2013.
- [17] J. Z. Leibo, V. Zambaldi, M. Lanctot, J. Marecki, and T. Graepel, “Multi-agent Reinforcement Learning in Sequential Social Dilemmas,” in *Proceedings of the 16th International Conference on Autonomous Agents and Multiagent Systems (AA-MAS 2017)*, (Sao Paulo, Brazil), 2017.
- [18] J. Perolat, J. Z. Leibo, V. Zambaldi, C. Beattie, K. Tuyls, and T. Graepel, “A multi-agent reinforcement learning model of common-pool resource appropriation,” in *Advances in Neural Information Processing Systems (NIPS)*, (Long Beach, CA), 2017.
- [19] J. N. Foerster, R. Y. Chen, M. Al-Shedivat, S. Whiteson, P. Abbeel, and I. Mordatch, “Learning with opponent-learning awareness,” *arXiv preprint arXiv:1709.04326*, 2017.
- [20] A. Lerer and A. Peysakhovich, “Maintaining cooperation in complex social dilemmas using deep reinforcement learning,” *arXiv preprint arXiv:1707.01068*, 2017.
- [21] A. Peysakhovich and A. Lerer, “Prosocial learning agents solve generalized stag hunts better than selfish ones,” *arXiv preprint arXiv:1709.02865*, 2017.
- [22] A. Peysakhovich and A. Lerer, “Consequentialist conditional cooperation in social dilemmas with imperfect information,” *CoRR*, vol. abs/1710.06975, 2017.
- [23] M. Kleiman-Weiner, M. K. Ho, J. L. Austerweil, M. L. Littman, and J. B. Tenenbaum, “Coordinate to cooperate or compete: Abstract goals and joint intentions in social interaction,” in *CogSci*, 2016.

- [24] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, and T. Graepel, “Value-decomposition networks for cooperative multi-agent learning,” *CoRR*, vol. abs/1706.05296, 2017.
- [25] J. J. Rousseau, *Discourse on the Origin of Inequality*. Marc-Michel Rey, 1755.
- [26] H. L. A. Hart, “Are there any natural rights?,” *The Philosophical Review*, vol. 64, no. 2, pp. 175–191, 1955.
- [27] J. Rawls, “Justice as fairness,” *The philosophical review*, vol. 67, no. 2, pp. 164–194, 1958.
- [28] G. Klosko, “The principle of fairness and political obligation,” *Ethics*, vol. 97, no. 2, pp. 353–362, 1987.
- [29] B. S. Frey and I. Bohnet, “Institutions Affect Fairness: Experimental Investigations,” *Journal of Institutional and Theoretical Economics*, vol. 151, pp. 286–303, June 1995.
- [30] C. Bicchieri and A. Chavez, “Behaving as expected: Public information and fairness norms,” *J. Behav. Decis. Making*, vol. 23, pp. 161–178, Apr. 2010.
- [31] J. Henrich, J. Ensminger, R. McElreath, A. Barr, C. Barrett, *et al.*, “Markets, Religion, Community Size, and the Evolution of Fairness and Punishment,” *Science*, vol. 327, pp. 1480–1484, Mar. 2010.
- [32] R. Gibbons, *A primer in game theory*. Harvester Wheatsheaf, 1992.
- [33] C. Eckel and H. Gintis, “Blaming the messenger: Notes on the current state of experimental economics,” *Journal of Economic Behavior & Organization*, vol. 73, no. 1, pp. 109–119, 2010.
- [34] G. Charness and M. Rabin, “Understanding social preferences with simple tests,” *The Quarterly Journal of Economics*, vol. 117, no. 3, pp. 817–869, 2002.
- [35] D. Engelmann and M. Strobel, “Inequality aversion, efficiency, and maximin preferences in simple distribution experiments,” *American economic review*, vol. 94, no. 4, pp. 857–869, 2004.
- [36] C. Camerer, *Behavioral Game Theory: Experiments in Strategic Interaction*. 01 2011.
- [37] G. R. Grice, “The relation of secondary reinforcement to delayed reward in visual discrimination learning,” *Journal of Experimental Psychology*, vol. 38, no. 1, pp. 1–16, 1948.
- [38] M. J. Kearns and S. P. Singh, “Bias-variance error bounds for temporal difference updates,” in *Proceedings of the Thirteenth Annual Conference on Computational Learning Theory, COLT ’00*, (San Francisco, CA, USA), pp. 142–147, Morgan Kaufmann Publishers Inc., 2000.
- [39] B. F. Skinner, *The Behavior of Organisms; An Experimental Analysis*. D. Appleton-Century Company, 1938.
- [40] L. S. Shapley, “Stochastic Games,” *In Proc. of the National Academy of Sciences of the United States of America*, 1953.
- [41] M. L. Littman, “Markov games as a framework for multi-agent reinforcement learning,” in *Proceedings of the 11th International Conference on Machine Learning (ICML)*, pp. 157–163, 1994.
- [42] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pp. 1928–1937, 2016.
- [43] W. E. Walsh, R. Das, G. Tesauro, and J. O. Kephart, “Analyzing complex strategic interactions in multi-agent systems,” in *AAAI-02 Workshop on Game-Theoretic and Decision-Theoretic Agents*, pp. 109–118, 2002.
- [44] M. P. Wellman, “Methods for empirical game-theoretic analysis,” in *Proceedings of the national conference on artificial intelligence*, vol. 21, p. 1552, Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2006.
- [45] T. C. Schelling, “Hockey helmets, concealed weapons, and daylight saving: A study of binary choices with externalities,” *The Journal of Conflict Resolution*, vol. 17, no. 3, pp. 381–428, 1973.
- [46] P. Kollock, “Social dilemmas: The anatomy of cooperation,” *Annual review of sociology*, vol. 24, no. 1, pp. 183–214, 1998.

- [47] P. Oliver, “Rewards and punishments as selective incentives for collective action: theoretical investigations,” *American journal of sociology*, vol. 85, no. 6, pp. 1356–1375, 1980.
- [48] Ö. Güreker, B. Irlenbusch, and B. Rockenbach, “The competitive advantage of sanctioning institutions,” *Science*, vol. 312, no. 5770, pp. 108–111, 2006.
- [49] N. Chentanez, A. G. Barto, and S. P. Singh, “Intrinsically motivated reinforcement learning,” in *Advances in neural information processing systems*, pp. 1281–1288, 2005.
- [50] G. F. Loewenstein, L. Thompson, and M. H. Bazerman, “Social utility and decision making in interpersonal contexts,” *Journal of Personality and Social psychology*, vol. 57, no. 3, p. 426, 1989.
- [51] C. Bellemare, S. Kröger, and A. Van Soest, “Measuring inequity aversion in a heterogeneous population using experimental decisions and subjective probabilities,” *Econometrica*, vol. 76, no. 4, pp. 815–839, 2008.
- [52] E. I. Hoppe and P. W. Schmitz, “Contracting under incomplete information and social preferences: An experimental study,” *The Review of Economic Studies*, vol. 80, no. 4, pp. 1516–1544, 2013.
- [53] K. Verbeeck, J. Parent, and A. Nowé, “Homo equalis reinforcement learning agents for load balancing,” in *Innovative Concepts for Agent-Based Systems, First International Workshop on Radical Agent Concepts, WRAC 2002, McLean, VA, USA, January 16-18, 2002, Revised Papers*, pp. 81–91, 2002.
- [54] S. de Jong and K. Tuyls, “Human-inspired computational fairness,” *Autonomous Agents and Multi-Agent Systems*, vol. 22, no. 1, pp. 103–126, 2011.
- [55] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 1st ed., 1998.
- [56] T. Yamagishi, “The provision of a sanctioning system as a public good,” vol. 51, pp. 110–116, 07 1986.
- [57] E. Fehr and S. Gächter, “Cooperation and punishment in public goods experiments,” *American Economic Review*, vol. 90, pp. 980–994, September 2000.
- [58] R. P. Blake, K. Mcauliffe, J. Corbit, T. Callaghan, O. Barry, A. Bowie, L. Kleutsch, K. Kramer, E. Ross, H. Vongsachang, R. Wrangham, and F. Warneken, “The ontogeny of fairness in seven societies,” vol. 528, 11 2015.
- [59] S. F. Brosnan and F. B. M. de Waal, “Evolution of responses to (un)fairness,” *Science (New York, N.Y.)*, vol. 346, no. 6207, p. 1251776, 2014.
- [60] G. Chalkiadakis and C. Boutilier, “Coordination in multiagent reinforcement learning: a bayesian approach,” in *The Second International Joint Conference on Autonomous Agents & Multiagent Systems, AAMAS 2003, July 14-18, 2003, Melbourne, Victoria, Australia, Proceedings*, pp. 709–716, 2003.
- [61] C. Gini, *Variabilità e mutabilità: contributo allo studio delle distribuzioni e delle relazioni statistiche*. No. pt. 1 in Studi economico-giuridici pubblicati per cura della facoltà di Giurisprudenza della R. Università di Cagliari, Tipogr. di P. Cuppini, 1912.
- [62] C. F. Camerer and E. Fehr, “Measuring social norms and preferences using experimental games: A guide for social scientists,” *Foundations of human sociality: Economic experiments and ethnographic evidence from fifteen small-scale societies*, vol. 97, pp. 55–95, 2004.