

Evolution of cooperation among individuals with limited payoff memory

Christian Hilbe, Nikoleta E. Glynatsi, Alex McAvoy

Abstract

1 Introduction

One of the most important applications of evolutionary game theory is the evolution of cooperation. Why is it that some individuals choose to help others (increasing their payoff) at the expense of decreasing one's own payoff? During the past decade, the literature on mechanisms that allow for cooperation, even though it's seen to be at odd has been fruitful. One such mechanism is repetition; the so called direct reciprocity. Theoretical researchers have been using evolutionary models to understand how direct reciprocity allows cooperation to evolve and which strategies are important for sustaining cooperation.

Evolutionary game theory does not require individuals to be rational, instead they adapt strategies based on mutation and exploration. Strategies with high fitness are more likely to spread, either because the individuals who adopt them have more offspring (fitness-based processes), or they are imitated more often (pairwise comparison processes) [1]. The fitness of strategies, and subsequently of the individuals, is not constant. Instead it depends on the composition of the population. Individuals interact with other members of the population according to their strategies. Their yielded payoffs, according to pre defined games, are translated into fitness.

It is commonly assumed that fitness assimilates to the *expected payoff*, which the mean payoff an individual achieved over the different types in the population after multiple intercalations. Expected payoffs imply that individuals have a perfect a memory. In order to estimate expected payoffs individuals must be able to recall several of their interactions with each player in the population. However, when modeling how they make decisions in each round they assumed to have very limited memory. To be precise, most work on this models focuses on naive subjects who can only choose from a restricted set of strategies [4], or who do not remember anything beyond the outcome of the very last round [5], with a few notable exceptions [2, 3].

This creates a curious inconsistency, which raises the following question: how robust is our understanding of cooperation? We are not the first to question the assumptions of models when estimating fitness. In [6]

relax the assumption that selection occurs much more slowly than the interaction between individuals. They results shows that rapid selection affects evolutionary dynamics in such a dramatic way that for some games it even changes the stability of equilibria.

We consider two cases. Initially, we start by considering two extreme scenarios. The first is the classical scenario of the expected payoffs and the alternative scenario where individuals update their strategies only based on the very last payoff they obtained. We refer to these as the *stochastic* payoffs of individuals. We do this for the donation games but also the important cases of symmetric 2×2 games. These include the prisoner's dilemma, the snowdrift game, the stag hunt game and the harmony game.

In the later sections we allow individuals to use more memory. More specifically, individuals update their strategies by considering up to the last two rounds, whilst interacting with up to two different players.

2 Model Setup

In order to account for the difference in the robustness of cooperation among individuals based on their payoff memory, we consider a population of N players, where N is even, and where mutations are sufficiently rare. At any point in time there are at most two different strategies present in the population. A *resident* strategy and a *mutant* strategy. Suppose there are $N - k$ players who use the resident strategy whereas k players use the mutant strategy. Each step of the evolutionary process consists of two stages, a game stage and an updating stage.

In the game stage, each player is randomly matched with some other player in the population to interact in a number of turns. The number of turns is not fixed, on the contrast we consider that another interaction has a probability δ to continue following each turn. At each turn the players can choose to either cooperate (C) or to defect (D). The payoffs in a given round depend on both player's decisions and are given by:

$$U = \begin{pmatrix} R & S \\ T & P \end{pmatrix}$$

We assume herein that individuals at most make use of simple *reactive strategies* make decisions in each round. Reactive strategy are a set of memory-one strategies that only take into account the previous action of the opponent. Reactive strategies can be written explicitly as a vector $\in \mathbb{R}_3$. More specifically a reactive strategy s is given by $s = (y, p, q)$ where y is the probability that the strategy opens with a cooperation and p, q are the probabilities that the strategy cooperates given that the opponent cooperated and defected equivalently.

In the updating stage, two players are randomly drawn from the population, a 'learner' and a 'role model'. Given that the learner's payoff $u_L \in \mathcal{U}$ and that the role model's payoff $u_{RL} \in \mathcal{U}$, we assume the learner adopts the role model's strategy based on the Fermi distribution function. The relative influence of the payoffs on the adopt the strategy of the other is controlled by an external parameter. The so called intensity of selection,

$\beta \geq 0$.

This basic evolutionary step is repeated until either the mutant strategy goes extinct, or until it fixes in the population, in which case the mutant strategy becomes the new resident strategy. After either outcome, we introduce a new mutant strategy, uniformly chosen from all reactive strategies at random, and we set the number of mutants to $k = 1$. This process of mutation and fixation/extinction is then iterated many times.

In this work we explore the effect on the payoffs in the updating stage. The details of the process described in this section can be found in appendix A.

3 Analysis of Stochastic Payoffs in the Donation game

We first explore the effect of stochastic payoffs on the cooperative behaviour by considering the extreme case that an individual's fitness is based on their last interaction with one other individual. We compare this to the expected payoffs. In this section we consider the donation game. Each player can cooperate by providing a benefit b to the other player at their cost c , with $0 < c < b$. Then, $T = b, R = b - c, S = -c, P = 0$, and matrix (2) is given by:

$$\begin{pmatrix} b - c & -c \\ b & 0 \end{pmatrix} \quad (1)$$

Figure 1 shows simulation results for the described process of section 2. Figure 1 depicts the evolving conditional cooperation probabilities p and q . Assuming that the discount factor δ is comparably high, the opening move y is a transient effect and has no effect on the outcome. The left panels correspond to the standard scenario considered in the previous literature. It considers players who use expected payoffs to update their strategies. The right panel shows the scenario considered herein, in which players update their strategies based on their last round's payoff. The top panels assume a benefit of 3, whereas the bottom panels assume that the benefit is 10.

The figure suggests that when updating is based on expected payoffs, players tend to be more generous. The q -values are higher on average which suggests that individuals tend to cooperate more after they are at the receiving end of a defection. In addition, for both cases of b , individuals tend to be more cooperative. The average cooperation rate is strictly higher for expected payoffs. This difference is statistically important, and the effect is even more obvious when the benefit is higher. More specifically, for $b = 10$ the average cooperation rate drops from 97% to 51%.

Figure 2 further explores the effect of the benefit on the cooperation rate. In all cases, the stochastic payoffs evolution tends to reduce the evolving cooperation rate. For the stochastic payoffs, benefit appears to not make a difference once $b > 5$, and on average the evolving cooperation rate is at 50%. The effect of benefit on the expected payoffs smooths out after $b > 6$, and the expected payoffs estimating the evolving rate to be at 90%.

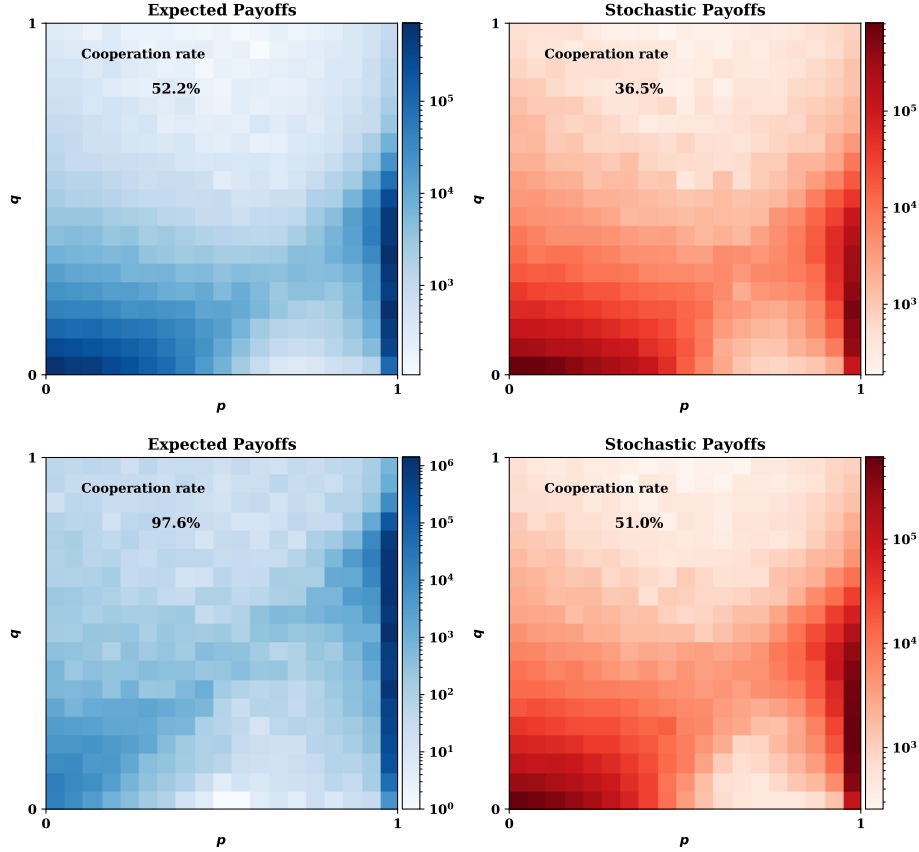


Figure 1: Evolutionary dynamics under expected payoffs and stochastic payoffs. We have run two simulations of the evolutionary process described in Section ?? for $T = 10^7$ time steps. For each time step, we have recorded the current resident population (y, p, q) . Since simulations are run for a relatively high continuation probability of $\delta = 0.999$, we do not report the players' initial cooperation probability y . The graphs show how often the resident population chooses each combination (p, q) of conditional cooperation probabilities in the subsequent rounds. (A) If players update based on their expected payoffs, the resident population typically applies a strategy for which $p \approx 1$ and $q \leq 1 - c/b = 0.9$. The cooperation rate within the resident population (averaged over all games and over all time steps) is close to 100%. (B) When players update their strategies based on their realized payoffs in the last round, there are two different predominant behaviors. The resident population either consists of defectors (with $p \approx q \approx 0$) or of conditional cooperators. In the latter case, the maximum level of q consistent with stable cooperation is somewhat smaller compared to the expected-payoff setting, $q < 0.5$. Also the resulting cooperation rate is smaller. On average, players cooperate roughly in half of all rounds. Parameters: $N = 100$, $b = 3$, $c = 1$, $\beta = 1$, $\delta = 0.999$.

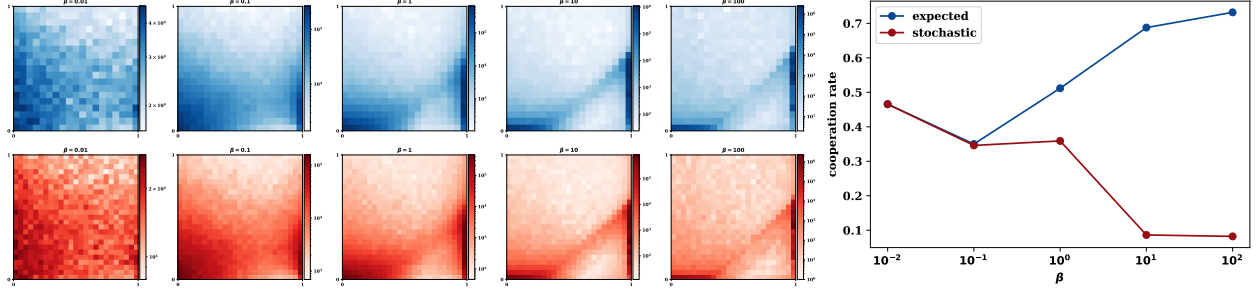


Figure 2: The evolution of cooperation for different benefit values. Here, we vary the benefit of defection b . In all cases, stochastic payoff evaluation tends to reduce the evolving cooperation rates. Unless explicitly varied, the parameters of the simulation are $N = 100$, $b = 3$, $c = 1$, $\beta = 1$, $\delta = 0.99$. Simulations are run for $T = 5 \times 10^6$ time steps for each parameter combination.

Figure 3 illustrates the results for different runs of the evolutionary process where we vary the strength of selection. In the case of a very small $\beta < \beta = 10^{-1}$ the process is almost random, as the effect of the payoffs is very small. Once payoff begin to matter the difference is once again evident, with the expected payoffs always overestimating the evolved cooperation rate.

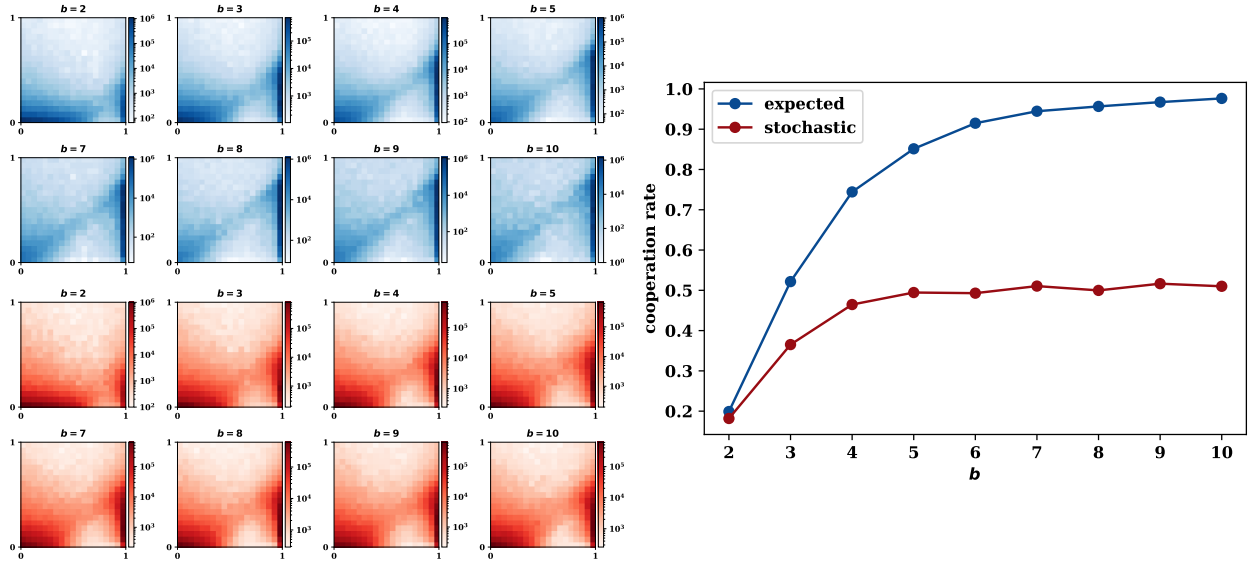


Figure 3: The evolution of cooperation for different selection strength values. Here, we vary the selection strength β . In all cases, stochastic payoff evaluation tends to reduce the evolving cooperation rates. Unless explicitly varied, the parameters of the simulation are $N = 100$, $b = 3$, $c = 1$, $\beta = 1$, $\delta = 0.99$. Simulations are run for $T = 5 \times 10^6$ time steps for each parameter combination.

4 Expected and stochastic payoffs in 2×2 games

A Model Setup

In order to account for the difference in the robustness of cooperation among individuals based on their payoff memory, we consider a population of N players, where N is even, and where mutations are sufficiently rare. At any point in time there are at most two different strategies present in the population. A *resident* strategy and a *mutant* strategy. Suppose there are $N - k$ players who use the resident strategy whereas k players use the mutant strategy. Each step of the evolutionary process consists of two stages, a game stage and an updating stage.

In the game stage, each player is randomly matched with some other player in the population to interact in a number of instances of the game $\begin{pmatrix} R & S \\ T & P \end{pmatrix}$. The number of interactions is not fixed, on the contrast we consider that another interaction has a probability δ to continue following each turn. We assume herein that individuals at most make use of simple *reactive strategies* make decisions in each round. Reactive strategies are a set of memory-one strategies that only take into account the previous action of the opponent. Reactive strategies can be written explicitly as a vector $\in \mathbb{R}_3$. More specifically a reactive strategy s is given by $s = (y, p, q)$ where y is the probability that the strategy opens with a cooperation and p, q are the probabilities that the strategy cooperates given that the opponent cooperated and defected equivalently.

In the updating stage, two players are randomly drawn from the population, a ‘learner’ and a ‘role model’. Given that the learner’s payoff $u_L \in \mathcal{U}$ and that the role model’s payoff $u_{RL} \in \mathcal{U}$, we assume the learner adopts the role model’s strategy with probability

$$\rho(u_L, u_{RL}) = \frac{1}{1 + \exp[-\beta(u_L - u_{RL})]}. \quad (2)$$

where $\beta \geq 0$ corresponds to relative influence of the payoffs on the adopt the strategy of the other is controlled by an external parameter, the so called intensity of selection.

We iterate this basic evolutionary step until either the mutant strategy goes extinct, or until it fixes in the population (in which case the mutant strategy becomes the new resident strategy). After either outcome, we introduce a new mutant strategy (uniformly chosen from all reactive strategies at random), and we set the number of mutants to $k = 1$. This process of mutation and fixation/extinction is then iterated many times. The fixation probability of the mutant strategy then takes the standard form [7],

$$\varphi = \frac{1}{1 + \sum_{i=1}^{N-1} \prod_k^i \frac{\lambda_k^-}{\lambda_k^+}}. \quad (3)$$

where λ_k^-, λ_k^+ are the probabilities that the number of mutants decreases and increases respectively.

We compare this process for stochastic payoff evaluation with the analogous process where players update their strategies with respect to their *expected* payoffs.

A.1 Expected Payoffs

Consider two players with reactive strategies $s_1 = (y_1, p_1, q_1)$ and $s_2 = (y_2, p_2, q_2)$ who interact in a repeated prisoner's dilemma with continuation probability δ , the probability that are in each of the four possible states in the last round is given by the stationary distribution $\langle \mathbf{v}$ of the transition matrix M .

$$M = \begin{bmatrix} p_1 p_2 & p_1 (1 - p_2) & p_2 (1 - p_1) & (1 - p_1) (1 - p_2) \\ p_2 q_1 & q_1 (1 - p_2) & p_2 (1 - q_1) & (1 - p_2) (1 - q_1) \\ p_1 q_2 & p_1 (1 - q_2) & q_2 (1 - p_1) & (1 - p_1) (1 - q_2) \\ q_1 q_2 & q_1 (1 - q_2) & q_2 (1 - q_1) & (1 - q_1) (1 - q_2) \end{bmatrix}. \quad (4)$$

The long run steady state probability vector $\langle \mathbf{v}$ is the solution to $\langle \mathbf{v} M = \langle \mathbf{v}$, can be combined with the payoff u to give the payoffs for each player. More specifically,

$$\langle \mathbf{v}(s_1, s_2), \mathbf{u} \rangle$$

is the payoff that s_1 achieves when interacting with s_2 . In the evolutionary process given k individuals with mutant strategy s_M and $N - k$ individuals with the resident strategy s_R the expected payoffs are,

$$\begin{aligned} \pi_R &= \frac{N - k - 1}{N - 1} \cdot \langle \mathbf{v}(s_R, s_R), \mathbf{u} \rangle + \frac{k}{N - 1} \cdot \langle \mathbf{v}(s_R, s_M), \mathbf{u} \rangle, \\ \pi_M &= \frac{N - k}{N - 1} \cdot \langle \mathbf{v}(s_M, s_R), \mathbf{u} \rangle + \frac{k - 1}{N - 1} \cdot \langle \mathbf{v}(s_M, s_M), \mathbf{u} \rangle. \end{aligned} \quad (5)$$

In the limit of no discounting, $\delta \rightarrow 1$, this process based on expected payoffs has been considered in [8].

A.2 Stochastic Payoffs

Initially, we consider the case that an individual's fitness is based on their last interaction with one other individual. There only four possible outcomes for the last round, those are CC, CD, DC, DD . Consider two players with reactive strategies $s_1 = (y_1, p_1, q_1)$ and $s_2 = (y_2, p_2, q_2)$ who interact in a repeated prisoner's dilemma with continuation probability δ , the probability that are in each of the four possible states in the last round is given by:

$$\mathbf{v}(s_1, s_2) = \left(\mathbf{v}_R(s_1, s_2), \mathbf{v}_S(s_1, s_2), \mathbf{v}_T(s_1, s_2), \mathbf{v}_P(s_1, s_2) \right). \quad (6)$$

where,

$$\begin{aligned}
\mathbf{v}_R(S_1, S_2) &= (1-\delta) \frac{y_1 y_2}{1-\delta^2 r_1 r_2} + \delta \frac{\left(q_1 + r_1((1-\delta)y_2 + \delta q_2) \right) \left(q_2 + r_2((1-\delta)y_1 + \delta q_1) \right)}{(1-\delta r_1 r_2)(1-\delta^2 r_1 r_2)}, \\
\mathbf{v}_S(S_1, S_2) &= (1-\delta) \frac{y_1 \bar{y}_2}{1-\delta^2 r_1 r_2} + \delta \frac{\left(q_1 + r_1((1-\delta)y_2 + \delta q_2) \right) \left(\bar{q}_2 + \bar{r}_2((1-\delta)y_1 + \delta p_1) \right)}{(1-\delta r_1 r_2)(1-\delta^2 r_1 r_2)}, \\
\mathbf{v}_T(S_1, S_2) &= (1-\delta) \frac{\bar{y}_1 y_2}{1-\delta^2 r_1 r_2} + \delta \frac{\left(\bar{q}_1 + \bar{r}_1((1-\delta)y_2 + \delta p_2) \right) \left(q_2 + r_2((1-\delta)y_1 + \delta q_1) \right)}{(1-\delta r_1 r_2)(1-\delta^2 r_1 r_2)}, \\
\mathbf{v}_P(S_1, S_2) &= (1-\delta) \frac{\bar{y}_1 \bar{y}_2}{1-\delta^2 r_1 r_2} + \delta \frac{\left(\bar{q}_1 + \bar{r}_1((1-\delta)y_2 + \delta p_2) \right) \left(\bar{q}_2 + \bar{r}_2((1-\delta)y_1 + \delta p_1) \right)}{(1-\delta r_1 r_2)(1-\delta^2 r_1 r_2)}.
\end{aligned} \tag{7}$$

Proof. Assume a repeated prisoner's dilemma between two reactive strategies. Given the continuation probability δ , probability that the game ends in the after the first round $(1 - \delta)$ and the expected distribution of the four outcomes in the very first round is \mathbf{v}_0 defined as. Following the first round the, the outcome of the next rounds with a probability δ is M such that,

...

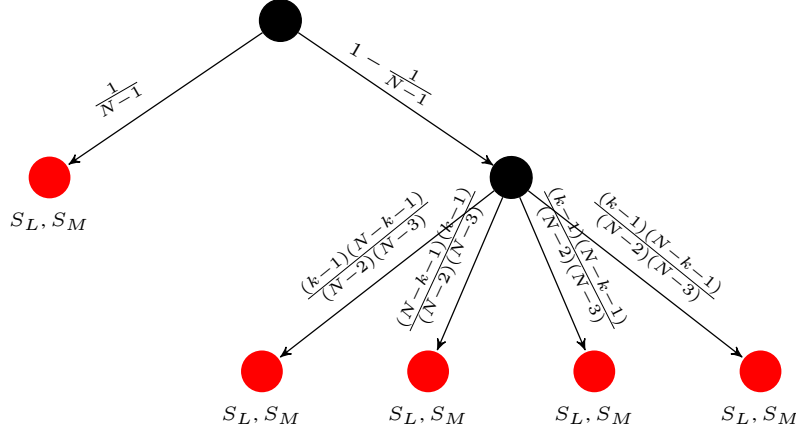
It can shown that, $(1-\delta)\mathbf{v}_0(I_4 - \delta M)^{-1}$ and with some algebraic manipulation we derive to Equation 7. \square

Equation 7 is the probability vector that players s_1, s_2 are in each of the possible states of the last round. Given the population N with k mutants in the last round there only possible combinations of interactions are:

A.3 Fixation probabilities under stochastic payoff evaluation

Given that $N - k$ players use the resident strategy $S_1 = (y_1, p_1, q_1)$ and that the remaining k players use the mutant strategy $S_2 = (y_2, p_2, q_2)$, the probability that the number of mutants increases by one in one step of the evolutionary process can be written as

$$\lambda_k^+ = \frac{N-k}{N} \cdot \frac{k}{N} \cdot \sum_{u_1, u_2 \in \mathcal{U}} x(u_1, u_2) \cdot \rho(u_1, u_2). \tag{8}$$



In this expression, $(N-k)/N$ is the probability that the randomly chosen learner is a resident, and k/N is the probability that the role model is a mutant. The sum corresponds to the total probability that the learner adopts the role model's strategy over all possible payoffs u_1 and u_2 that the two player may have received in their respective last rounds. We use $x(u_1, u_2)$ to denote the probability that the randomly chosen resident obtained a payoff of u_1 in the last round of his respective game, and that the mutant obtained a payoff of u_2 . Given that the payoffs are u_1 and u_2 , the imitation probability is then given by $\rho(u_1, u_2)$, as specified by Eq. (2). The probability that the respective payoffs of the players are given by u_1 and u_2 can be calculated as

$$\begin{aligned}
 x(u_1, u_2) = & \frac{1}{N-1} \cdot v_{u_1}(S_1, S_2) \cdot 1_{(u_1, u_2) \in \mathcal{U}_F^2} \\
 & + \left(1 - \frac{1}{N-1}\right) \left[\frac{k-1}{N-2} \frac{k-2}{N-3} v_{u_1}(S_1, S_2) v_{u_2}(S_2, S_2) + \frac{k-1}{N-2} \frac{N-k-1}{N-3} v_{u_1}(S_1, S_2) v_{u_2}(S_2, S_1) \right. \\
 & \quad \left. + \frac{N-k-1}{N-2} \frac{k-1}{N-3} v_{u_1}(S_1, S_1) v_{u_2}(S_2, S_2) + \frac{N-k-1}{N-2} \frac{N-k-2}{N-3} v_{u_1}(S_1, S_1) v_{u_2}(S_2, S_1) \right]. \tag{9}
 \end{aligned}$$

The first term on the right side corresponds to the case that the learner and the role model happened to be matched during the game stage, which happens with probability $1/(N-1)$. In that case, we note that only those payoff pairs can occur that are feasible in a direct interaction, $(u_1, u_2) \in \mathcal{U}_F^2 := \{(R, R), (S, T), (T, S), (P, P)\}$, as represented by the respective indicator function. Otherwise, if the learner and the role model did not interact directly, we need to distinguish four different cases, depending on whether the learner was matched with a resident or a mutant, and depending on whether the role model was matched with a resident or a mutant.

References

- [1] Bin Wu, Benedikt Bauer, Tobias Galla, and Arne Traulsen. Fitness-based models and pairwise comparison models of evolutionary games are typically different—even in unstructured populations. *New Journal of Physics*, 17(2):023043, 2015.
- [2] Ch Hauert and Heinz Georg Schuster. Effects of increasing the number of players and memory size in the iterated prisoner’s dilemma: a numerical approach. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 264(1381):513–519, 1997.
- [3] Alexander J Stewart and Joshua B Plotkin. Small groups and long memories promote cooperation. *Scientific reports*, 6(1):1–11, 2016.
- [4] Martin A Nowak and Karl Sigmund. Tit for tat in heterogeneous populations. *Nature*, 355(6357):250–253, 1992.
- [5] Seung Ki Baek, Hyeong-Chai Jeong, Christian Hilbe, and Martin A Nowak. Comparing reactive and memory-one strategies of direct reciprocity. *Scientific reports*, 6(1):1–13, 2016.
- [6] Carlos P. Roca, José A. Cuesta, and Angel Sánchez. Time scales in evolutionary dynamics. *Phys. Rev. Lett.*, 97:158701, Oct 2006.
- [7] Martin A Nowak, Akira Sasaki, Christine Taylor, and Drew Fudenberg. Emergence of cooperation and evolutionary stability in finite populations. *Nature*, 428(6983):646–650, 2004.
- [8] Lorens A Imhof and Martin A Nowak. Stochastic evolutionary dynamics of direct reciprocity. *Proceedings of the Royal Society B: Biological Sciences*, 277(1680):463–468, 2010.