

The Evolution of Cooperation in the Centipede Game with Finite Populations*

Rory Smead^{†‡}

The partial cooperation displayed by subjects in the Centipede Game deviates radically from the predictions of traditional game theory. Even standard, infinite population, evolutionary settings have failed to provide an explanation for this behavior. However, recent work in finite population evolutionary models has shown that such settings can produce radically different results from the standard models. This paper examines the evolution of partial cooperation in finite populations. The results reveal a new possible explanation that is not open to the standard models and gives us reason to be cautious when employing these otherwise helpful idealizations.

1. Introduction. The *Centipede Game*, first introduced by Rosenthal (1981), provides a setting where rational self-interest conflicts with socially optimal cooperative behavior. Two players have to divide an increasing sum of money; to do this, they take turns choosing between taking the larger share of the money for themselves or passing the choice to the other player; every time a player passes on her opportunity the sum increases. The game ends when a player takes the larger amount of money or after a set number of turns, in which case the money is divided in a prearranged manner. It is set up such that if a self-interested player believes the game will end on the next round, she should take the money and end it on the current round. Several of the most popular solution concepts in traditional game theory tell us that we should expect to see absolutely no

*Received January 2007; revised November 2007.

[†]To contact the author, please write to: Department of Logic and Philosophy of Science, University of California, 3151 Social Science Plaza A, Irvine, CA 92697-5100; e-mail: rsmead@uci.edu.

[‡]I would like to thank Brian Skyrms, Kevin Zollman, Michael McBride, two anonymous referees, and the members of the Social Dynamics Seminar at UCI for their helpful feedback on this paper. Generous financial support was provided by the School of Social Sciences at UCI.

Philosophy of Science 75 (April 2008) pp. 157–177. 0031-8248/2008/7502-0003\$10.00
Copyright 2008 by the Philosophy of Science Association. All rights reserved.

cooperation (passing) in the Centipede Game. However, experimental results deviate wildly from these expectations; subjects seem invariably to exhibit partial (though, not total) cooperation.

Rosenthal (1981) initially argued that we should treat Centipede-like situations as pair of single person decision problems. And, from a decision-theoretic standpoint, if each player in the Centipede Game has certain views about the propensity of the other to deviate from the game-theoretic predictions, then it will be rational to pass in early rounds (96–97). But why should each expect the other to deviate from the game-theoretic solution? We could augment Rosenthal's analysis with the hypothesis that humans simply have a natural propensity to cooperate, at least partially, in situations like the Centipede Game. This leaves us with a puzzle: how can we explain this propensity?

For such a question, it is natural to turn to evolutionary game theory. However, the standard evolutionary models do not provide an answer, yielding the same solution as the traditional game-theoretic concepts and so do not aid us in explaining the existence of partial cooperation. However, the standard evolutionary setting assumes an infinite population, and recent studies into finite population evolutionary dynamics have revealed some important differences between these finite cases and their infinite counterpart. For instance, Imhof, Fudenberg, and Nowak (2005) show that stochastic evolution in finite populations need not select for a strictly dominant strategy. And Taylor et al. (2004) show that this evolutionary dynamic can even select *against* such strategies in 2×2 games.¹

I will examine the evolution of partial cooperation in finite populations using the Moran process (Moran 1962) with frequency-dependent fitness. Through the use of computer simulations, I find that, in finite populations, evolutionary paths often 'get stuck' in states where the population is exhibiting partial cooperation rather than the state of no cooperation, as seen in the infinite setting. This result illustrates a new potential explanation for the evolution of partial cooperation as well as providing a striking example of how finite population evolutionary models can differ radically from the traditional models. The frequently used and mathematically useful idealizations of the infinite population setting can lead us away from possible evolutionary explanations of partial cooperation. More generally, this suggests that one ought to be cautious when relying on such idealizations in evolutionary models.

1. The effect Taylor et al. (2004) notice is that since individuals in finite populations do not interact with themselves, this creates a small degree of anti-correlation between strategies causing the cross-strategy payoffs to be slightly more important to fitness than same-strategy payoffs. This means that in some settings, finite populations can favor spiteful (but dominated) strategies.

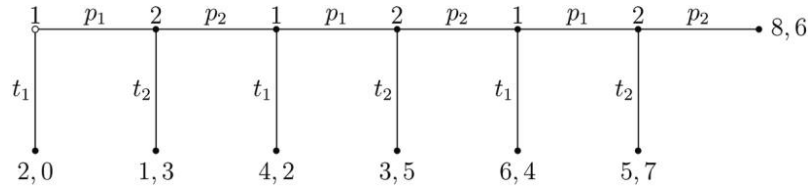


Figure 1. A six-stage Centipede Game.

I will begin by setting out two games of partial cooperation: the Centipede Game and the Quasi-Centipede Game. In Section 3, I will examine these games in the standard evolutionary setting using the traditional replicator dynamic. And, in Section 4, I will explore a model of evolution in finite populations and the effect of this finite setting on partial cooperation.

2. The Centipede and Quasi-Centipede Games. As mentioned above, the Centipede Game is a simple case where self-interest can interfere with more efficient, cooperative behavior. Imagine two players sitting at a table with a sum of money α between them, say \$2. They take turns choosing either to pass (p) or to take a significant portion, leaving the rest to the other player (t). This continues for a fixed, and known, number of rounds or until someone takes the money. If a player chooses to pass, the sum increases by \$2; if a player chooses to take the money, that player receives $\beta = 1/2(\alpha) + \$1$ and the other player receives the remaining money ($\alpha - \beta$). If no player chooses to take the money, it is divided as if the player who would have gone next takes the money. The payoffs are such that if the game ends immediately after a player passes, then that player would have been better off taking the money and ending the game.

In the classical example, the game continues for one hundred rounds; hence the name. However, any game with the basic structure and payoff ordering is a Centipede Game. Figure 1 shows a six-stage version of this game.

This can be easily generalized. Let the strategy sets be $S_1 = \{1, 3, 5, \dots, N + 1\}$ and $S_2 = \{2, 4, 6, \dots, N + 1\}$, where N is the length of the Centipede Game.² And, let $s_i \in S_i$ denote the strategy chosen by player i ; this number corresponds to the round on which the player decides to defect (if $s_i = N + 1$, then player i passes at every stage). Suppose

2. This is actually a reduced strategy set relative to the traditional Centipede Game. Each strategy should include what the player would choose at every choice point, whether or not that choice point is actually reached. However, for our purposes in this paper this reduced strategy set will suffice.

$s_i = s_j - 1$; then the utility function for player i is such that $\pi_i(s_i + 2, s_j) < \pi_i(s_i, s_j)$. In other words, if my opponent is going to defect on the next round, I am better off defecting now rather than passing. The payoff functions from the above six-stage Centipede Game can be generalized as

$$\pi_1(s_1, s_2) = \begin{cases} s_1 + 1 & \text{if } s_1 \leq s_2 \\ s_2 - 1 & \text{if } s_1 > s_2 \end{cases},$$

$$\pi_2(s_2, s_1) = \begin{cases} s_1 - 1 & \text{if } s_1 \leq s_2 \\ s_2 + 1 & \text{if } s_1 > s_2 \end{cases}.$$

The Centipede Game is often used in game theory courses to introduce the concepts of *backward induction* and the *iterated elimination of weakly dominated strategies*. Suppose player 2 were to find herself at the last stage in the Centipede Game above; naturally she would prefer to take the money. But, given that player 2 would take the money at the last stage, player 1 should take the money on the second-to-last stage and given player 1's preference, player 2 should defect on the third-to-last stage. Iterating this reasoning leads to the conclusion that player 1 should take the money on the first stage. The unique *subgame-perfect equilibrium* of this game is for each player to opt out at any chance she gets. Every *Nash equilibrium* (Nash 1950) in this game will share the feature of first-stage defection with the subgame-perfect equilibrium.³ It is this uniquely 'rational' (and unintuitive) solution that traditional game-theoretic wisdom tells us to expect.

Given that the payoffs for a few rounds of cooperation (passing) are so much better than just taking the money immediately, it seems unlikely that people would actually conform to the 'rational' solution. The counterintuitive nature of the backwards induction solution has led to debate over its rationality.⁴ Furthermore, in experimental settings, subjects do tend to show some degree of cooperation. McKelvey and Palfrey (1992) studied the behavior of subjects in both four-stage and six-stage Centipede Games at varying stakes. They find that subjects display a substantial amount of cooperation, typically making it more than halfway down the

3. There are actually a large number of Nash equilibria in the Centipede Game. Any pair of (mixed) strategies where player 1 defects immediately and player 2 would choose to defect on her first choice frequently enough to make player 1's immediate defection a best-response will qualify as a Nash equilibrium. However, the only pure-strategy Nash Equilibria is where $s_1 = 1$ and $s_2 = 2$.

4. Some papers in this debate include Binmore 1994, 1996 and Aumann 1995, 1996.

centipede before one takes the money (808).⁵ They also found that higher stakes slightly lowered the degree of cooperation, but by a relatively small amount.⁶

The fact that people regularly deviate from the predictions of traditional game theory has led many scholars to look for alternative explanations. For instance, McKelvey and Palfrey show that if we assume that some individuals are altruists (unconditional passers), then it may be within the scope of 'rational' self-interest to pass in early rounds. McKelvey and Palfrey also look at models with errors in action or belief and show that partial cooperation could also arise in these settings. Nagel and Tang (1998) look to various models of learning with the hope of providing some insight. They find that a model of simple reinforcement learning does fairly well in predicting partial cooperation. These explanations may certainly be part of the answer, but there may be alternative (or supplementary) explanations which do not rely on errors, preexistence of altruism, or specific learning mechanisms. Organisms of all kinds have solved cooperation problems, of which Centipede Games are a specific kind, and it would be worthwhile to investigate how partial cooperation could evolve in even the most basic settings. This is the aim of the models presented below.

For the sake of illustration, the model that will be investigated below will largely focus on a simultaneous move and symmetric game similar to the Centipede Game, which I will call a 'Quasi-Centipede Game'.⁷ This game begins with a sum of money (say \$4) between the players and at each round they simultaneously declare to 'take' or 'pass'. If both players choose 'take' then they split the money and the game is over. If either chooses to 'pass' the amount of money will increase (at \$2 for each 'pass')

5. In the Centipede Game studied by McKelvey and Palfrey, the payoffs grew exponentially with the pot, doubling every time a player chose to pass rather than by a fixed amount, as in the game in Figure 1. In the McKelvey and Palfrey game there is more incentive to cooperate, but choosing to 'take' also gives a fixed proportion of the pot so there is also greater incentive to deviate as well. The results that will be presented in this paper will also hold for these exponentially increasing Centipedes and they turn out to be, in some sense, easier for cooperation (see the Appendix). Here we will focus on the Centipede Game with steady growth in payoff because it presents the (slightly) harder case for partial cooperation.

6. For a survey of various experimental results on dominance-solvable games like the Centipede Game see Camerer 2003, 218–236.

7. This modified game simplifies the model and, as will be shown, the qualitative results will hold for the standard Centipede Game as well, meaning the results here are robust across different games of similar form. Also, experimental results of subjects playing a game that is somewhat similar to the one above closely resemble those of the standard Centipede Game. See Van Huyck, Wildenthal, and Battalio (2002) for these results; the term 'Quasi-Centipede' is used by Camerer (2003).

TABLE 1. A THREE-STAGE QUASI-CENTIPEDE GAME

	$s_2 = 0$	$s_2 = 1$	$s_2 = 2$
$s_1 = 0$	2, 2	5, 1	5, 1
$s_1 = 1$	1, 5	4, 4	7, 3
$s_1 = 2$	1, 5	3, 7	6, 6

action). If one chooses ‘take’ and the other ‘pass’, the taker receives the majority (\$4 more than the other player), leaving the rest to the other and ending the game. If they both choose ‘pass’ then the sum increases and they play again. This continues until one or both chooses ‘take’ or for a finite and known number of rounds, at which point the money is split equally. As with the standard Centipede Game, we will consider a restricted strategy set where the player simply decides which round to choose ‘take’. Table 1 shows a three-stage Quasi-Centipede Game; the possible strategies are simply the number of rounds that each player chooses to ‘pass’ before they ‘take’.

The payoff function in this game can easily be generalized to the N -stage game, here is one example:

$$\pi_1(s_1, s_2) = \begin{cases} 2s_1 + 5 & \text{if } s_1 < s_2 \\ 2s_1 + 2 & \text{if } s_1 = s_2, \\ 2s_2 + 1 & \text{if } s_1 > s_2 \end{cases}$$

where $0 \leq s_i < N$ is player i 's strategy; the payoff function for player two is similar. This payoff function will be used in the model examined below.⁸ To see the relationship between the Centipede Game and Quasi-Centipede Game, it is helpful to see the Quasi-Centipede Game in the extensive form (Figure 2). Here s_i represents the number of times the player chooses p_i or ‘pass’. For the backward induction, begin at the end of the game and notice that player 2 should choose t_2 no matter what player 1 chose, and given that player 2 will choose t_2 , player 1 should choose t_1 and so on.⁹

Once again, the unique solution for this game is to take the money immediately, and this solution can be delivered by an iterated elimination of weakly dominated strategies. But, in some sense, the force of this solution is even stronger than in the Centipede Game since there is a unique Nash Equilibrium at $(s_1 = 0, s_2 = 0)$. The standard Centipede Game permits some degree of cooperation from player 2 in the equilibria of the game, since when player 1 chooses ‘take’ right away, player 2 never

8. This formulation of the Quasi-Centipede Game is similar to another game known as the Traveler’s Dilemma, which has similar solutions and similar empirical results (Basu 1994; Capra et al. 1999).

9. The Quasi-Centipede Game is somewhat similar to a standard Centipede Game where the players do not know their position.

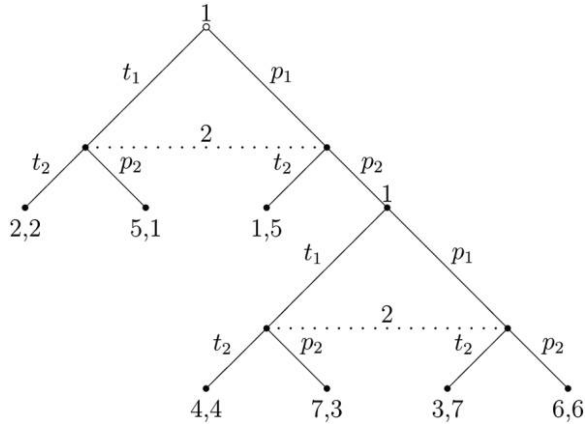


Figure 2. A three-stage Quasi-Centipede Game in extensive form.

has the opportunity to act. In the Quasi-Centipede with the restricted strategy set above, *both* players are strictly worse when they deviate from equilibrium behavior. Furthermore, if models of the Quasi-Centipede Game yield similar results to the Centipede Game, this can be taken as an indication that the results are robust across different games that have similar structures.

3. Evolution in the Centipede Game. For the evolutionary setting, we imagine an infinite population of players that are ‘hard-wired’ to cooperate for a certain number of rounds in the Centipede Game. Individuals are paired randomly to play the game and receive payoffs which function as the ‘fitness’ of using a particular strategy k . Let $0 \leq x_s \leq 1$ represent the frequency of strategy s in the population ($\sum_{s \in S} x_s = 1$) and $X = (s_0, s_1, \dots)$ represent the state of the population. The fitness of type $s' \in S$ (the set of individuals who use strategy s') is calculated by

$$f_{s'}(X) = \sum_{s \in S} \pi(s', s)x_s.$$

Let θ_X be the average fitness of the population in state X ; the distribution of strategies evolves according to the standard replicator dynamic.¹⁰ The differential equation governing the change in frequency of type s is

$$\dot{x}_s = x_s(f_s(X) - \theta_X).$$

In words, types that have higher fitness increase proportionally to their

10. For details on this dynamic, see Hofbauer and Sigmund (1998).

fitness at the expense of those with lower fitness. The replicator dynamic is most directly interpreted as a model of biological evolution, but it can also be interpreted in terms of imitation rather than reproduction and hence be used to model a form of cultural evolution. For this cultural interpretation, ‘fitness’ would refer to a propensity for a strategy to be imitated rather than for an individual of that type to reproduce.

The replicator dynamic does not generally support the elimination of weakly dominated strategies, a fact that can lead to interesting results. For instance, Gale, Binmore, and Samuelson (1995) and Skyrms (1996) show that weakly dominated ‘fair’ strategies can survive in evolutionary models of the Ultimatum Game. Thus, there may be some hope for promising results in the Centipede Game using the replicator dynamic. However, in both the Centipede Game and Quasi-Centipede Game the replicator dynamic yields a single result: a population of first-round defectors.

In the case of the three-stage Quasi-Centipede Game above, it is easy to see that there will be a unique evolutionary outcome from the replicator dynamic: everyone plays 0. If we suppose the population consists almost entirely of type 2, it will be invaded by individuals of type 1 and if the population consists of almost entirely type 1 then it will be invaded by individuals of type 0.¹¹ The same result holds for arbitrarily large Quasi-Centipede Games. Figure 3 shows the global dynamics for a three-stage Quasi-Centipede Game. It is interesting to note that the dynamics do not necessarily carry the population directly to the unique equilibrium. Instead, the dynamic sometimes seems carry the population through the stages of an elimination of dominated strategies, first to a predominance of type 2, then to type 1, and finally to the evolutionary equilibrium at all type 0.

As seen in Figure 3, some paths of evolution in this three-stage game come close to a nonequilibrium edge of the simplex. The longer the N -stage Quasi-Centipede Game, the more likely it is to have an evolutionary path that gets close to a nonequilibrium edge of the simplex. The reason for this is that there is no direct route from type $s_i = k$ to type $s_i = k - 2$ (for any $N \geq k \geq 2$) and so, to evolve the unique solution as mentioned above, the population must evolve through the stages of backwards induction.

The evolutionary results for the standard Centipede Game are similar

11. This is under the assumption that every type is represented in the initial population. The replicator dynamics will then only carry some types to extinction in the limit, a point that will be important when we turn to finite populations.

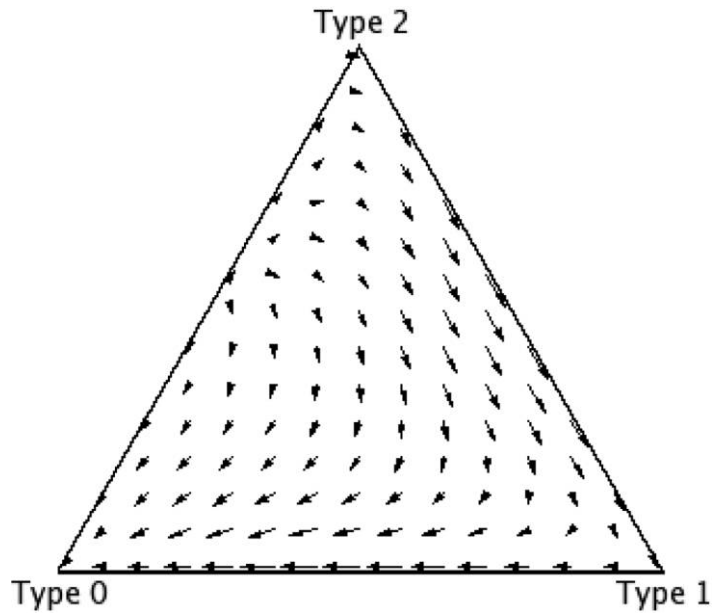


Figure 3. Simplex of a three-stage Quasi-Centipede Game.

to the Quasi-Centipede Game.¹² In a sufficiently mixed population, types that show a high-degree of cooperation initially outperform those that do not and hence, increase in frequency. Then, those types that exhibit slightly less cooperation grow at the expense of the more cooperative type; and this process repeats until the population consists only of first-round defectors. In general, stable states of the replicator dynamic must also be Nash equilibria and all Nash equilibria in the Centipede Game have the property of first-round defection.¹³ Ponti (2000) shows that first round defection in the Centipede Game will be the result in a large class of dynamics that includes the replicator dynamic; proving that these dynamics will eventually carry the population to a state that is 'outcome equivalent' to the subgame perfect equilibrium (total first-round defection).

12. The possible types have to be adjusted for the Centipede Game. Specifically, an individual's strategy must include what the individual should do if in the role of player 1 and what to do if in the role of player 2.

13. Some weakly dominated strategies will survive in small amounts. In a population of first-round defectors, there is no difference in payoff between strategies with regard to the action when in the role of player 2. Thus, some strategies which are cooperative when in the role of player 2 may survive in stable states, but there is always unanimous first-round defection when in the role of player 1.

Ponti also notes that populations often exhibit ‘unlearning’ of the solution (increasing in cooperation) before converging on first-round defection and he shows that as the length of the Centipede increases, the more likely populations are to exhibit ‘unlearning’.

At this point we can conclude that the standard evolutionary models (the replicator dynamic) cannot alone account for the evolution of partial cooperation. Even mindless replicators can evolve the backwards-induction solution to the Centipede Game. However, if we include some perturbation in the dynamics (which could be interpreted as either mutation in the case of biological evolution or experimentation in the case of cultural evolution), we no longer see the first-round defection result. In addition to his results mentioned above, Ponti also shows that with enough perturbation in the dynamic, a population can exhibit cycles between higher and lower degrees of cooperation.

The cyclic dynamical picture given by mutation in the Centipede Game provides one potential explanation for the existence and evolution of partial cooperation. In the following section I will examine a model that will provide another possible evolutionary explanation that does not rely on mutation.

4. Partial Cooperation in Finite Populations. In recent years, Taylor et al. (2004) have developed a stochastic model of evolution in finite populations by incorporating frequency-dependent fitness into the Moran process (Moran 1962).¹⁴ Taylor et al. have shown that this process can deliver radically different results than the standard replicator dynamic.¹⁵ I will use the Moran process to model the evolution of partial cooperation in the Centipede and Quasi-Centipede Games. The frequency-dependent Moran process provides a finite population evolutionary setting similar in nature to the standard evolutionary setting, which allows us to examine, in particular, the assumption of infinite populations. I find that without any mutation (or a very small amount), resulting populations will regularly exhibit partial cooperation.

On this model, the population consists of M individuals. Let $x'_s \leq M$ denote the number of individuals of type s at time t ; then $\sum_{s \in S} x'_s = M$.

14. This process has also been studied in a series of papers by Fudenberg et al. (2004, 2006) and has been applied to the repeated prisoners dilemma by Nowak et al. (2005).

15. Others in evolutionary game theory have also explored different finite population models with interesting philosophical results. For instance, Grim, Mar, and Denis (1998), Alexander and Skyrms (1999), Alexander (2000), Skyrms (2004), Vanderschraaf and Alexander (2005), and Zollman (2005) have explored the evolution of finite population models in local interaction settings. However, the model that will be investigated here is an unstructured population with a different dynamic.

And let $X_t = (x_{s_1}^t, x_{s_2}^t, \dots)$ represent the state of the population at time t . The fitness to a given type s' at time t in the population is

$$f_{s'}^t(X_t) = \pi(s', s')(x_{s'}^t - 1) + \sum_{s \neq s'} \pi(s', s)x_s^t.$$

The evolution is stochastic. At each step, one individual is chosen to reproduce with a probability proportional to her fitness and produces one identical offspring. Then, another individual is selected randomly to die. The probability for reproduction of type s' at time t is

$$R_{s'}^t(X_t) = \frac{x_{s'}^t f_{s'}^t}{\sum_s x_s^t f_s^t},$$

and the probability for a death of a type s is simply x_s^t/M . We can also talk of selection favoring certain strategies at certain states of the population, we can say that selection favors x_s in state X_t if $R_s^t(X_t) > R_{s'}^t(X_t)$ for all $s' \neq s$.

As the population size becomes very large ($M \rightarrow \infty$) and time is appropriately adjusted, this process is similar, though not identical, to the standard replicator dynamic (Traulsen, Claussen, and Hauert 2005).¹⁶ But for finite populations, a population governed by the Moran process is effectively a very large finite-state Markov chain with as many absorbing states as there are types in the population. Once $x_s^t = M$ for some s , that type has completely overtaken the population and becomes the only type that can reproduce ($R_s^t = 1$); hence the population stays in this state. Also, if the frequency x_s^t of some type s drops to $x_s^t = 0$, then it cannot increase because $R_s^t = 0$, and so extinction of a type is possible on this dynamic. As with the replicator dynamic, the Moran process can be interpreted as biological evolution in terms of reproduction (with a population at the carrying capacity of the environment) or as cultural evolution in terms of imitative behavior. I will intentionally remain ambiguous with respect to the interpretation since the target here is traditional idealizations and generic potential explanations rather than an explanation for a specific case of partial cooperation. Each ‘generation’ of the Moran process is

16. Traulsen et al. (2005) give a detailed analysis of the limiting behavior of the Moran process and how it compares to the replicator dynamic, showing that in the limit the Moran process is identical to the adjusted replicator dynamic from Maynard Smith 1982.

just a single change in an individual, either through death and reproduction or through imitation.¹⁷

We can now ask: what happens to partial cooperation in this evolutionary setting of finite populations? Moving to finite populations complicates the mathematics significantly, making analytical results difficult to achieve. Consequently, the one very effective way to investigate these models is through computer simulations. The advantage is a more realistic evolutionary setting, with respect to population sizes, but we are limited to examining the numerical results of specific models rather than the generality of mathematical proofs.

4.1. The Quasi-Centipede in Finite Populations. To answer this question, I will first focus on the Quasi-Centipede Game and then turn to the Centipede Game. I will assume that the Quasi-Centipede Game being played is the same as the generalized one mentioned above and will employ the use of simulations. Just as in the traditional evolutionary setting, we will not necessarily get a direct trajectory to the evolutionary equilibrium; instead, an appropriately mixed population will take a round-about path similar to some shown in Figure 3. However, under the Moran process we see several differences. First, the dynamic is discrete-time rather than continuous. Second, the step-direction is stochastic, only approximating the force of selection on average. Third, as a result of these differences, it is almost always possible for a type to go extinct in this setting.

The fact that extinction is possible and that the dynamic is stochastic can now cause ‘breaks’ in the evolutionary chain that follows the iterated elimination of dominated strategies. This can result in the population being absorbed into a state where every member is playing a partially cooperative strategy. Recall Figure 3; in this setting we can imagine an evolutionary path close to the Type 2–Type 1 edge of the simplex where Type 0 happens to go extinct.¹⁸ In this case, the population will most likely be absorbed into a state where every individual is playing Type 1 (since Type 1 outperforms Type 2 in this population). If extinction of strategies in the Quasi-Centipede Game is possible, then a population can ‘get stuck’ along its evolutionary path and end up playing a partially cooperative strategy. But, how likely is such an outcome?

17. There are other finite population dynamics that may be of interest as well; for the purposes of this paper I focus only on the Moran process because it has already generated results in other areas of game theory (Fudenberg et al. 2004, 2006; Nowak et al. 2005) and it retains many features of the replicator dynamic.

18. For instance, there could be only one Type 0 in a population of Type 1’s and 2’s; Type 1 could be chosen to reproduce and the lone Type 0 could be chosen to die leaving $x_0 = 0$.

The answer to this question depends on many factors: the size of the population, the length of the Quasi-Centipede, and the initial state of the population. I will begin with simulations in an example case. Suppose $M = 1,000$ and suppose that each individual is given a random strategy, with equal probability for each strategy (equal initial weights); for a 10-stage Centipede Game, simulations show an overwhelming likelihood of resulting populations ending in an absorbing state with a substantial degree of cooperation. In 1,000 trials, there was not a single case of the unique equilibrium seen in the traditional setting and most ended in partial cooperation over the half-way point in the Centipede in an average of just over 150,000 generations (iterations of the Moran process).

One may object that this method of assigning starting strategies is overly biased toward the center of the population's state-space. This can be remedied by selecting a random point on the simplex of strategies and assigning the weights accordingly (random initial weights).¹⁹ This modification causes more variance in terms of which absorbing state is reached, but we still observe partial cooperation with very high frequency. Figure 4 gives a summary of the results for 1,000 trials under both styles of random initial conditions. In each case, we can see that it is very likely (probability greater than 0.99) that an evolving population will become absorbed into a state where every member exhibits partial cooperation.

As mentioned above, there are other key variables that affect the resulting degree of cooperation in a population: the length of the Centipede and the size of the population.²⁰ To gauge the degree of cooperation across different games we can assign it a number in $[0, 1]$ which is determined by the number of 'passes' by members of a population divided by the number of possible 'passes' in the game; call this the 'level of cooperation'. If we increase the population size M , then the average level of cooperation seen in the resulting population decreases. The reason is that the step-size ($1/M$) of the evolution becomes proportionally smaller as M increases, and so the population is less likely to get extremely close to the edge of the simplex.

If, on the other hand, the size N of the Quasi-Centipede is increased, then we see the average level of cooperation in the resulting population increase. The reason is that, as N increases, the number of places extinction of strategies could occur increases. Since the evolution roughly follows

19. This is done by assigning each strategy s a random weight w_s on the interval $[0, 1]$ and assigning an individual strategy s with probability $-\ln(w_s)/[\sum_{s' \in S} -\ln(w_{s'})]$.

20. The specific values of the payoffs also make a difference in the outcome. We can, while preserving the ordinal ranking of payoffs, make things harder or easier for cooperation to result by making the payoff for defecting larger or smaller. However, the qualitative results are robust across such changes.

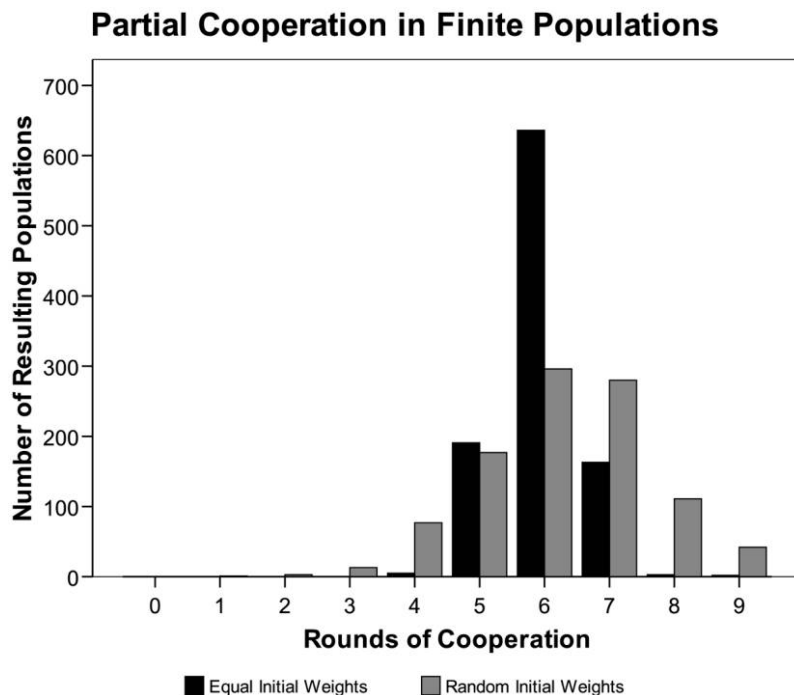


Figure 4. The results of 1,000 simulations on a 10-stage Quasi-Centipede Game with a population of 1,000 evolving by the frequency-dependent Moran process.

the stages of iterated elimination of weakly dominated strategies, there will be $N - 1$ edges of the simplex that correspond to steps in the elimination of dominated strategies. And, each time the evolution comes near an edge of the simplex, there is a chance of strategy extinction, and hence, a chance of becoming absorbed into a partially cooperative state.

Figure 5 shows the average level of cooperation (normalized on $[0,1]$ where 0 corresponds to no cooperation and 1 to complete cooperation) as a function of the population on a logarithmic scale for 5-, 7-, 10-, and 14-stage Quasi-Centipede Games. These simulations show that the larger the population, the more the dynamics resembles the infinite case, as expected. And, the longer the Centipede the higher the expected degree of cooperation in the resulting population.²¹ Given these results, it is

21. With very small Quasi-Centipede Games (such as $N = 3$) and the payoffs as in Table 1, populations resulting in partial cooperation are rarely observed. However, increasing the potential benefit for cooperation can increase the frequency of partial cooperation.

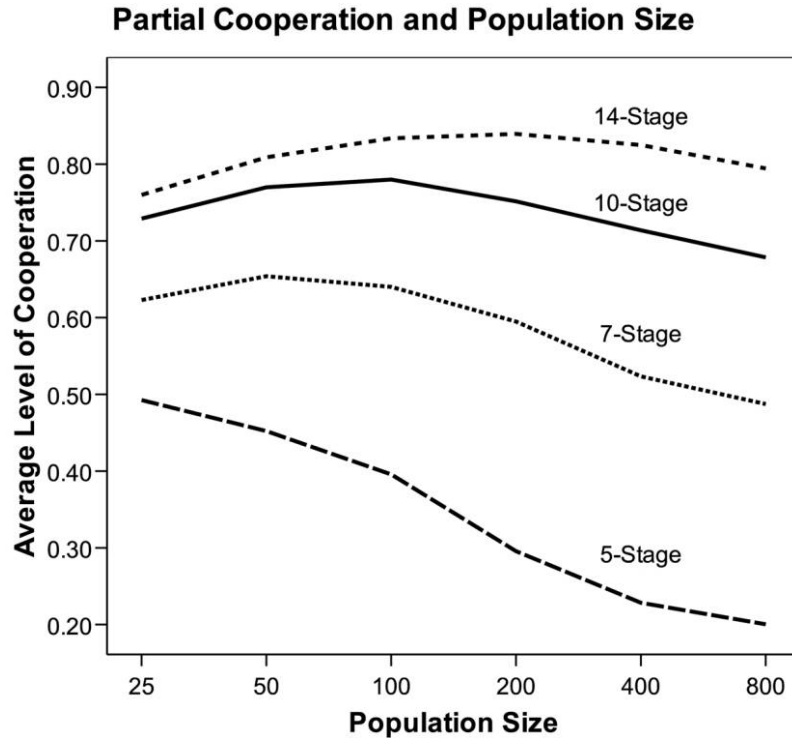


Figure 5. The average resulting levels of cooperation in 5-, 7-, 10-, and 14-stage Quasi-Centipede Games as a function of population size.

reasonable to conjecture that for any finite population size M there will be some N -stage Centipede Game (with payoffs as specified above) that will result in partial cooperation with a probability arbitrarily close to 1.

4.2. The Centipede Game in Finite Populations. Thus far, I have focused on the Quasi-Centipede Game. Simulations on the Centipede Game using the Moran process show that the same qualitative results hold for this game. For these results, we imagine two populations playing the Centipede Game specified in Section 2, where members from one always take the role of player 1 and members from the other always take the role of player 2. The fitness for a type in a population is calculated by how well they do against the other population and the evolution within a population is governed by the Moran process.²²

22. The single population setting simply has more types, where an individual's type

In an eight-stage Centipede Game being played between two populations of $M = 1,000$, as in the case of the Quasi-Centipede, there were no resulting populations of first-round defectors in 1,000 trials.²³ And, the average level of cooperation was over the half-way point of the Centipede (the average node that resulting populations played ‘take’ was 6.12) with the average time to an absorbing state was just under 200,000 generations. We also find that the same general qualitative results hold as in the Quasi-Centipede Game, the average level of cooperation in resulting populations increases with the length of the Centipede and decreases as population size M increases. As before, these changes are due to the changes in likelihood of strategies going extinct.

Thus, in both the Centipede and Quasi-Centipede Games, the Moran process illuminates a new possibility for the evolution of partial cooperation. Since the evolution tends to follow the stages of elimination of dominated strategies, it is easy for certain strategies to become extinct. This can cause the population to become absorbed into a partially cooperative state, and these simulations have shown it is very likely to do so. Thus, the only reason certain levels of cooperation persist is because the ‘superior’ but slightly less cooperative strategies went extinct prior to arrival at the current state of the population.²⁴

Furthermore, given the nature of the standard infinite population replicator dynamic, as long as the initial population begins in the interior of the state-space (where every type is represented in the population), types will go extinct only in the limit of the process. Thus, it is impossible to observe outcomes in the Centipede or Quasi-Centipede Game similar to those seen with the frequency-dependent Moran process. The move to finite populations, in the case of the Centipede Game (and Quasi-Centipede Game), opens up new possibilities of explanation.

4.3. *Mutation.* Although the Moran process is stochastic, the models

must specify what to do in the role of player 1 and in the role of player 2. The results for the single population are similar to the two-population case. However, single populations often persist in awkward polymorphic mixes for prolonged periods of time before reaching an absorbing state. In these mixes, every member defects on some round n when in the role of player 1 but there is no consensus on which round $m > n$ to defect on when in the role of player 2. For this reason, only results of the two-population case are presented here.

23. The initial conditions for these simulations were determined by assigning each individual a random strategy with equal probability across strategies. If strategies were assigned with random weights, as above, the same qualitative results seen in the Quasi-Centipede Game hold: more variety in resulting populations is observed, but an overwhelming majority show high degrees partial cooperation.

24. It is worthwhile to note that we have been left with a purely explanatory model, which carries no important normative lessons.

examined so far have not included any mutation. Indeed, the lack of mutation plays a key role in finite populations resulting in partially cooperative states since this result relies on the extinction of certain strategies. We can introduce a mutation probability μ into the reproduction process such that, with probability μ , a new individual picks a random strategy (with equal weights to each strategy) instead of adopting the strategy of the reproduced type.

By introducing a mutation rate, the process becomes ergodic: over time the system can reach any state from any other state. With a very small mutation rate, simulations show the results presented above tend to break down as the number of generations gets very large. In this setting, the population settles into a normal partially-cooperative state but eventually just the right mutation(s) will occur to cause a slightly less cooperative type to take over. However, for these very small mutation rates, it may take a very long time to reach the immediate-defection state. For instance, simulations in a 10-stage Quasi-Centipede Game with a population of $M = 1,000$ and a mutation chance of 1-in-10,000, can take several million generations before reaching the immediate-defection state. Furthermore, if the mutation rate is high enough (say, 0.02 in the 10-stage, $M = 1,000$ case), then simulations show a finite-population analogue to Ponti's (2000) cycles of evolution seen in infinite populations with perturbation.²⁵ Figure 6 shows the average level of cooperation over time in such a cycling population through the first one million generations.

Given that the evolution is stochastic, and the level of cooperation graphed in Figure 6 is an average of the members of a mixed population, some cycles appear more well behaved than others. The increase in level of cooperation tends to be more dramatic and the decrease in cooperation more gradual. This is due to the fact that one population is always invaded by any number of individuals that are slightly less cooperative, whereas for an invasion of more cooperative individuals there has to be several mutants that are substantially more cooperative. Increases in cooperation require several mutants so they can benefit from one and other, and they need to be substantially more cooperative otherwise they would benefit the native members more than mutants of their own kind. Similar cycles are seen in the standard Centipede Game.

Evolutionary cycles of cooperation have been seen in other games studied with the Moran process in finite populations. Imhof et al. (2005) have examined the repeated prisoner's dilemma in this setting and have shown that mutation can cause cycles between cooperative, uncooperative, and discriminatory strategies. In the Centipede and Quasi-Centipede Games,

25. If we were to interpret the Moran process here as a form of imitation (cultural evolution), then a 2% experimentation rate may be very reasonable.

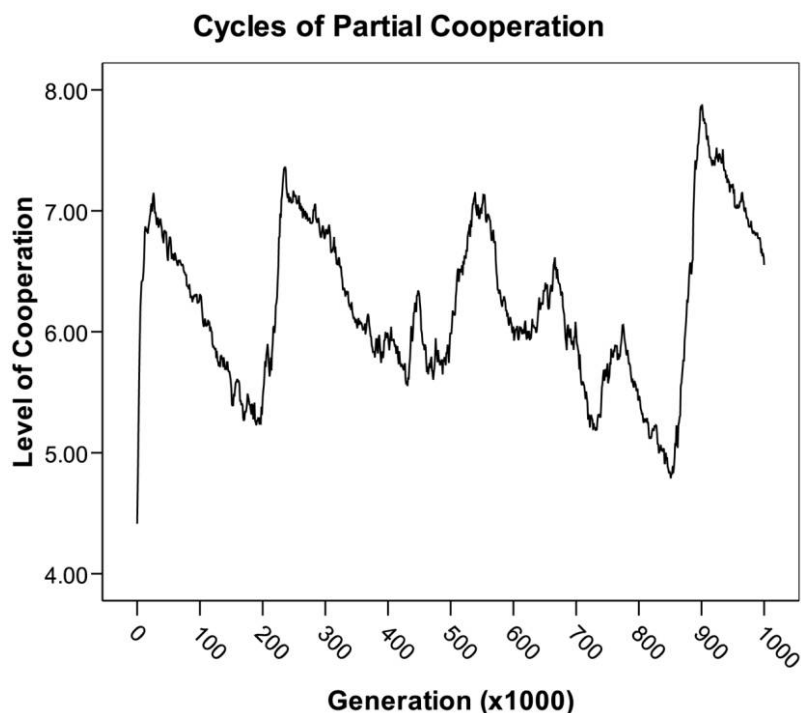


Figure 6. Cycles of partial cooperation in an evolving population of 1,000 individuals playing a 10-stage Quasi-Centipede Game with a mutation rate of $\mu = 0.02$.

cycles are possible with any mutation rate, but simulations show that as the mutation rate decreases, the regularity of cycles also decreases. There is a window of mutation that will tend to cause partial cooperation to break down and result in populations arriving at the immediately-defect solution (mimicking the infinite population setting without mutation). If $\mu = 0$ or is very small, then the population will get stuck in states of partial cooperation indefinitely or for very long periods of time. If μ is large, then we see frequent and indefinite cycles of partial cooperation in the population. Of course, the exact place and size of this window will vary with population size and the length of the Centipede.

5. Conclusion. Moving from infinite population models to finite population models can have a radical effect on the results of evolution. The possibility of extinction and the stochastic nature of finite population models opens up new explanations for behavior in certain games. We have seen that this is certainly the case with the Centipede and Quasi-Centipede

Games. In these games, the evolution roughly follows the stages of iterated elimination of weakly dominated strategies and in the standard infinite population models results in no cooperation. However, because of the nature of finite population evolutionary dynamics, certain strategies can go extinct, causing the population to ‘get stuck’ in partially cooperative states. The likelihood of establishing partial cooperation increases with the size of the game and decreases with the size of the population.

This result illustrates a new possible explanation for the evolution of partial cooperation: it may be that evolutionary paths, although leading in the direction noncooperative states, become absorbed into partially cooperative states due to the finite and stochastic nature of the evolutionary process. Although this explanation has not been shown to be more plausible than others, this result does provide a striking example of how evolutionary processes in finite settings can deviate significantly from those in the traditional infinite setting. The idealizations of the latter render such evolutionary outcomes as seen in finite cases impossible, and consequently, blind us to possible explanations for partial cooperation.

There are also several remaining questions that warrant further investigation. First, given that we see partial cooperation arise in both the Centipede and Quasi-Centipede, perhaps these methods could be applied to other games involving backwards induction. Second, how might different sorts of finite population dynamics affect the evolutionary picture in the setting of the Centipede Game or others? For instance, agent-based local interaction models may generate interesting results for the Centipede Game, and such models would enable further investigation of other idealizations in the traditional setting, in particular, random interaction within the population.

Appendix

The aforementioned results of the standard Centipede Game are presented in Table A1. The results are qualitatively similar to the Quasi-Centipede Game presented above, but the model has some additional subtleties due to the asymmetric nature of the game. Summarized below are 1,000 simulations of two 1,000-member populations playing an eight-legged standard Centipede Game with the same payoff structure as in Figure 1 (Regular Centipede) as well as an eight-legged Centipede Game with exponentially increasing payoffs (Exponential Centipede).²⁶ Each pair of populations

26. The payoffs for the exponential version are the same as in McKelvey and Palfrey 1992, beginning with a pot of 50, which doubles for each ‘pass’ and if a player chooses to ‘take’, then that player receives 4/5 of the pot while the other player receives 1/5.

gets absorbed into a single pair of strategies and we can examine where these populations first choose to ‘take’ in the course of the game (Round of First ‘Take’). The initial weights for strategies are given in parenthesis. The numerical values in the table represent the total number of populations that resulted in that level of cooperation (or round of defection) out of 1,000 simulated populations.

The move to an exponentially growing payoffs does not change the qualitative results of the Quasi-Centipede Game either. Table A2 summarizes the results for 1,000 simulations of a 1,000-member population playing a 10-stage Quasi-Centipede Game with exponentially increasing payoffs.

In both settings, the average level of cooperation increases (slightly) with the change to exponential payoffs and the behavior becomes more regular (less variance). It is in this sense that the case of exponential payoff increase partial cooperation becomes easier.

TABLE A1. RESULTS FOR THE STANDARD CENTIPEDE GAME.

Round of First “Take”	1	2	3	4	5	6	7	8	9
Reg. Centipede (Equal)	0	0	0	0	23	831	146	0	0
Reg. Centipede (Random)	2	1	10	44	170	347	309	98	19
Exp. Centipede (Equal)	0	0	0	0	0	356	644	0	0
Exp. Centipede (Random)	0	0	0	8	57	350	529	48	8

TABLE A2. RESULTS FOR AN EXPONENTIALLY INCREASING QUASI-CENTIPEDE GAME.

Rounds of Cooperation	0	1	2	3	4	5	6	7	8	9
Exp. Quasi (Equal)	0	0	0	0	0	0	844	156	0	0
Exp. Quasi (Random)	0	0	0	2	13	98	372	416	81	18

REFERENCES

- Alexander, J. McKenzie (2000), “Evolutionary Explanations of Distributive Justice”, *Philosophy of Science* 67: 490–516.
- Alexander, J. McKenzie, and Brian Skyrms (1999), “Bargaining with Neighbors: Is Justice Contagious?”, *Journal of Philosophy* 96: 588–598.
- Aumann, Robert J. (1995), “Backward Induction and Common Knowledge of Rationality”, *Games and Economic Behavior* 8: 6–19.
- (1996), “A Reply to Binmore”, *Games and Economic Behavior* 17: 138–146.
- Basu, Kaushik (1994), “The Traveler’s Dilemma: Paradoxes of Rationality in Game Theory”, *American Economic Review* 84: 391–395.
- Binmore, Kenneth (1994), “Rationality in the Centipede”, in Ronald Fagin (ed.), *Theoretical Aspects of Reasoning about Knowledge: Proceedings of the Fifth Conference (TARK 1994)*. San Francisco: Kaufmann, 150–159.
- (1996), “A Note on Backward Induction”, *Games and Economic Behavior* 17: 135–137.

- Camerer, Colin F. (2003), *Behavioral Game Theory: Experiments in Strategic Interaction*. Princeton, NJ: Princeton University Press.
- Capra, Monica C., Jacob K. Goeree, Rosario Gomez, and Charels A. Holt (1999), "Anomalous Behavior in a Traveler's Dilemma?", *American Economic Review* 89: 678–690.
- Fudenberg, Drew, Lorens A. Imhof, Martin A. Nowak, and Cristine Taylor (2004), "Stochastic Evolution as a Generalized Moran Process". Unpublished manuscript.
- Fudenberg, Drew, Martin A. Nowak, Cristine Taylor, and Lorens A. Imhof (2006), "Evolutionary Game Dynamics in Finite Populations with Strong Selection and Weak Mutation", *Theoretical Population Biology* 70: 352–363.
- Gale, John, Kenneth Binmore, and Larry Samuelson (1995), "Learning to Be Imperfect: The Ultimatum Game", *Games and Economic Behavior* 8: 56–90.
- Grim, Patrick, Gary Mar, and Paul St. Denis (1998), *The Philosophical Computer: Exploratory Essays in Philosophical Computer Modeling*. Cambridge, MA: MIT Press.
- Hofbauer, Josef, and Karl Sigmund (1998), *Evolutionary Games and Population Dynamics*. Cambridge: Cambridge University Press.
- Imhof, Lorens A., Drew Fudenberg, and Martin A. Nowak (2005), "Evolutionary Cycles of Cooperation and Defection", *PNAS* 102: 10797–10800.
- Maynard Smith, John (1982), *Evolution and the Theory of Games*. Cambridge: Cambridge University Press.
- McKelvey, Richard D., and Thomas R. Palfrey (1992), "An Experimental Study of the Centipede Game", *Econometrica* 60: 803–836.
- Moran, Patrick A. P. (1962), *The Statistical Processes of Evolutionary Theory*. Oxford: Clarendon.
- Nagel, Rosemarie, and Fang-Fang Tang (1998), "Experimental Results on the Centipede Game in Normal Form: An Investigation on Learning", *Journal of Mathematical Psychology* 42: 356–384.
- Nash, John (1950), "Equilibrium Points in N-Person Games", *Proceedings of the National Academy of Sciences* 36: 48–49.
- Nowak, Martin A., Akira Sasaki, Cristine Taylor, and Drew Fudenberg (2005), "Emergence of Cooperation and Evolutionary Stability in Finite Populations", *Nature* 428: 646–650.
- Ponti, Giovanni (2000), "Cycles of Learning in the Centipede Game", *Games and Economic Behavior* 30: 115–141.
- Rosenthal, Robert W. (1981), "Games of Perfect Information, Predatory Pricing, and the Chain Store", *Journal of Economic Theory* 25: 92–100.
- Skyrms, Brian (1996), *The Evolution of the Social Contract*. Cambridge: Cambridge University Press.
- (2004), *The Stag Hunt and the Evolution of Social Structure*. Cambridge: Cambridge University Press.
- Taylor, Cristine, Drew Fudenberg, Akira Sasaki, and Martin A. Nowak (2004), "Evolutionary Game Dynamics in Finite Populations", *Bulletin of Mathematical Biology* 66: 1621–1644.
- Traulsen, Arne, Jens Christian Claussen, and Christopher Hauert (2005), "Coevolutionary Dynamics: From Finite to Infinite Populations", *Physical Review Letters* 95: 238701.
- Vanderschraaf, Peter, and J. McKenzie Alexander (2005), "Follow the Leader: Local Interactions with Influence Neighborhoods", *Philosophy of Science* 72: 86–113.
- Van Huyck, John B., John M. Wildenthal, and Raymond C. Battalio (2002), "Tacit Cooperation, Strategic Uncertainty, and Coordination Failure: Evidence from Repeated Dominance Solvable Games", *Games and Economic Behavior* 38: 156–175.
- Zollman, Kevin J. S. (2005), "Talking to Neighbors: The Evolution of Regional Meaning", *Philosophy of Science* 72: 69–85.