



THE UNIVERSITY OF  
**SYDNEY**

Economics Working Paper Series

2015 - 5

Shared intentions: the evolution  
of collaboration

Jonathan Newton

December 2015

# Shared intentions: the evolution of collaboration

Jonathan Newton<sup>a,1</sup>

<sup>a</sup>*School of Economics, University of Sydney.*

---

## Abstract

The ability to form shared intentions and adjust one's choices in collaboration with others is a fundamental aspect of human nature. We discuss the forces that act for and against the evolution of this ability. In contrast to altruism and other non-fitness maximizing preferences, for large classes of games the ability to form shared intentions proliferates when rare without requiring group selection or assortativity in matching.

*Keywords:* Shared intentions, collaboration, evolution, game theory

**JEL Classification Numbers:** C73.

---

---

<sup>1</sup>Come rain or come shine, I can be reached at jonathan.newton@sydney.edu.au, telephone +61293514429. This work was completed while the author was supported by a Discovery Early Career Researcher Award funded by the Australian Research Council (Grant Number: 130101768), and originated in discussions with Sung-Ha Hwang about my work with Simon Angus. In addition to the support of these two, my work on this topic has benefited greatly from the comments and criticisms of Ingela Alger, Erik Mohlin, Herb Gintis, Natalie Gold, Toby Handfield, Atsushi Kajii, Kohei Kawamura, Christoph Kuzmics, Heinrich Nax, Ryoji Sawa, Peyton Young; as well as seminar audiences at University of Edinburgh, University of Oxford, University of Cologne, University of Bielefeld, Australian National University, University of Sydney, Monash University, University of Aizu, Kyoto University, Keio University, University of Tokyo and Hitotsubashi University.

*“Yet how much and how correctly would we think if we did not think, as it were, in community with others to whom we communicate our thoughts, and who communicate theirs with us!”*

– Immanuel Kant (1786)

## 1. Introduction

Humans are a collaborative species. We collaborate for good and for ill, motivated by love, hate, spite, envy, self-aggrandisement and the basic urges to feed and to reproduce. The understanding of collaboration and cooperation has long been a goal of economics. The current paper models the ability to collaborate as the ability to jointly optimize. That is, can *we* choose what is best for *us* rather than merely making decisions as individuals? It is shown that the ability to form *shared intentions* and take such joint decisions could have evolved amongst ancient populations who lacked the foresight and reasoning abilities of modern humans, and moreover, that this could happen even in circumstances hostile to the evolution of other behavioral types such as cooperators or altruists who might be expected to behave in ways which appear collaborative.

It has been argued in the philosophical literature that the intentions behind collective acts can be distinct from an aggregation of individual intentions (Bratman, 1992; Searle, 1990; Tuomela and Miller, 1988). This is “shared intentionality”, the idea that “we intend to do X” is distinct from “I intend to do X [because I think that she also intends to do X]”.<sup>2</sup> Through conversation, pointing and gesturing, or alternative forms of reasoning such as ‘team reasoning’ (Bacharach, 1999, 2006; Sugden, 2000) people form shared intentions. When combined with notions of optimization, shared intentions naturally give rise to collective agency. To see this, consider Alice and Bob who wish to take a

---

<sup>2</sup>There is disagreement amongst philosophers as to what extent shared intentions can be reduced to individual intentions. See also Butterfill (2012); Gilbert (1990); Gold and Sugden (2007); Velleman (1997). We take no position on this. Our results hold regardless of how agents form shared intentions.

drink together at one of two bars, *Grandma's* and *Stitch*. Both Alice and Bob prefer *Grandma's* to *Stitch*. Now imagine Alice stating “*I intend to go to Stitch because I think that Bob intends to go to Stitch.*” Such an intention is optimal from Alice’s perspective, given her beliefs about Bob’s intentions, regardless of the Pareto suboptimality of *Stitch* as a venue. Now, were Alice instead to state “*We intend to go to Stitch,*” then there exists a perfectly valid criticism: given that both Alice and Bob prefer *Grandma's* to *Stitch*, and neither has any incentive to deceive the other, it is irrational for them (as a plural entity) to hold such an intention. Economists will recognize this reasoning as similar to that underpinning concepts in game theory such as the Core (Gillies, 1959), Strong Equilibrium (Aumann, 1959), Coalition Proofness (Bernheim, Peleg and Whinston, 1987), Coalitional Rationalizability (Ambrus, 2009), Renegotiation Proofness (Farrell and Maskin, 1989) and Coalitional Stochastic Stability (Newton, 2012).

This paper demonstrates how conditions faced by paleolithic hunter-gatherer societies could have led to the evolution of the ability to collaboratively share intentions. On the one hand, the existence of problems that could be solved by collective action would have spurred the evolution of the ability to form shared intentions. On the other hand, those who could not participate in collaborative acts could sometimes free ride on the successes of others. This free riding would work against the evolution of the ability to share intentions. Note that the sharing of intentions and joint optimization is a *mutualistic* behavior: all participants gain from engaging in it. This does not prevent free riding, as third parties can obtain positive externalities from the collaboration of others, for example if Alice and Bob collaborate in hunting a buffalo, but Colm eats some of the leftovers. The mutualistic nature of jointly intentional behavior can be contrasted with altruistic behavior, in which one party sacrifices fitness for the benefit of another. It has been documented in the anthropology literature that much of the cooperation observed in hunter-gatherer societies is mutualistic. See Smith (2003) for a survey.

A consequence of the mutualistic nature of the sharing of intentions is that

such behavior can proliferate when rare. This is in stark contrast to cooperator types or altruists, who become extinct in similar circumstances. Furthermore, even when conditions are adverse to the evolution of the sharing of intentions, for example when there are many opportunities for free riding, some amount of sharing of intentions will persist in the population, a minority behavior that can then spread when conditions become favourable. Note that unlike models of the evolution of altruism and other non-fitness maximizing behaviors, neither repeated interaction (Trivers, 1971), nor kin-selection (Fisher, 1930; Hamilton, 1963), nor assortativity of interaction (Alger and Weibull, 2013; Eshel and Cavalli-Sforza, 1982; Wilson and Dugatkin, 1997), nor group selection (Bowles, 2006; Choi and Bowles, 2007; Haldane, 1932) is required for shared intentions to evolve.

It is hard to overstate the importance of shared intentions to human behavior. Recent work in developmental psychology has shown that from early childhood, human subjects display the ability and desire to engage in collaborative activities. This collaborative urge emerges prior to sophisticated logical inference and the ability to articulate hierarchical beliefs (Tomasello and Rakoczy, 2003, and citations therein). Moreover, the inclination towards collaborative behaviors is considerably weaker in non-human great apes (Tomasello and Carpenter, 2007; Tomasello and Herrmann, 2010).<sup>3</sup> This accumulated evidence has lent support to the hypothesis that human collaborative activity provided a niche in which a uniquely human cognition, replete with sophisticated modes of reasoning, could evolve. This is known as the shared intentionality hypothesis (Call, 2009) or the Vygotskian intelligence hypothesis (Moll and Tomasello, 2007; Tomasello, 2014; Vygotsky, 1980). The results of the current paper add to the plausibility of this hypothesis, as they show how even in populations of unsophisticated agents, collaborative behavior can evolve.

The author knows of only two other works that deal directly with the topic

---

<sup>3</sup>See also Wobber, Herrmann, Hare, Wrangham and Tomasello (2014); Tomasello, Carpenter, Call, Behne and Moll (2005) and the accompanying critical responses.

of the current paper<sup>4,5</sup>: [Bacharach \(2006, Chapter 3\)](#) and the study of [Angus and Newton \(2015\)](#). [Bacharach \(2006\)](#) gives a predominantly non-quantitative argument as to why a group selection mechanism would lead to collaborative ‘team reasoning’ in coordination problems and social dilemmas. However, in a simulations-based study of coordination games on networks, [Angus and Newton \(2015\)](#) show that group selection is far from sufficient for the evolution of collaboration, and that selective pressure *against* the sharing of intentions can arise at a group level due to the possibility of collaborative behavior slowing technological advance. The cited papers focus on multiple pairwise interactions for which payoffs are given by an underlying two player game. The current paper does not restrict itself to pairwise interaction and gives analytic results for a setting in which members of a population are randomly matched to play  $m$ -player games. In contrast to previous work, there is no group selection and selective pressure against the sharing of intentions arises from either (i) free riding on the positive externalities of collaboration by others, or (ii) negative externalities of collaboration on other potential collaborators. Finally, it is instructive to compare the evolution of shared intentions to the evolution of preferences (e.g. [Dekel, Ely and Yilankaya, 2007](#); [Güth and Kliemt, 1998](#); [Robson, 1996](#); [Samuelson, 2001](#)). In contrast to the evolution of preferences, the ability to collaboratively share intentions does not change individuals’ ranking of outcomes. Instead it makes new outcomes available to individuals when they update their strategies as part of a group. Any individual’s ranking of menu items does not change, but the variety of items on the menu becomes more appealing.

---

<sup>4</sup>We emphasize that we are considering the evolution of a trait - the ability to collaborate and share intentions, *not* the evolution of the play of any specific ‘cooperative’ action. Alice and Bob may intend to plan a surprise party for Colm, or to rob him of his possessions. Either way, Alice and Bob are collaborating, but to quite different ends.

<sup>5</sup>An alternative approach to understanding collaboration is that of [Gavrilets \(2014\)](#), who models collaborative ability as entering directly into production functions. Groups with high levels of collaborative ability produce more of a public good, giving an advantage in a group selection framework.

The paper is organized as follows. Section 2 gives the model and considers the evolution of shared intentions for a large class of games which includes threshold public goods games,  $m$ -player Prisoner’s Dilemmas, trust games, the centipede game, Nash demand games, Bertrand oligopoly, minimum effort games and, under an additional condition, finitely or infinitely repeated versions of all of the above. Section 3 analyzes the evolution of shared intentions when collaboration can exert negative externalities on others, such as when two people team up to steal from a third party. Section 4 considers a continuum of types distinguished by different probabilities of an individual of a given type being in a collaborative frame of mind. Section 5 compares and contrasts our results to those for altruism and other behavioral explanations for ‘cooperation’ found in the literature. Section 6 concludes. All proofs not in the main text are relegated to the appendix.

## 2. Model and analysis

We shall consider a population of individuals represented by the unit interval. Fitnesses will be determined when randomly formed groups of  $m$  individuals encounter problems. These problems could be opportunities to hunt large prey such as whales (Alvard, 2001; Alvard and Nolin, 2002), or the possibility that coordinated action could bring about a large haul of small prey, such as is the case with fishing (Sosis, Feldstein and Hill, 1998).

### 2.1. The game

Formally, we represent a problem faced by a group of  $m$  individuals by  $\Gamma$ , an  $m$ -player game with player set  $M = \{1, \dots, m\}$  and strategy sets  $S_i$ ,  $i \in M$ . Let  $s_i \in S_i$  and  $s = (s_1, s_2, \dots, s_m)$  be representative strategies and strategy profiles respectively. Let  $S = \times_{i \in M} S_i$  be the set of all strategy profiles. Let  $\pi_i(s)$  be the payoff of player  $i$  at strategy profile  $s$ . That is,  $\pi_i(\cdot) : S \rightarrow \mathbb{R}$ . Payoffs represent reproductive fitness. Let  $\underline{x} := (\times, \dots, \times)$  be the status quo strategy profile and a Nash equilibrium of the game, where, for every  $i \in M$ ,

	+ <sub>1</sub>	+ <sub>2</sub>	×
+ <sub>1</sub>	$b - c$	$-c$	$-c$
+ <sub>2</sub>	$-c$	$b - c$	$-c$
×	0	0	0

(i) Two stags

	+ <sub>1</sub>	+ <sub>2</sub>	×
+ <sub>1</sub>	$-c$	$b - c$	$-c$
+ <sub>2</sub>	$b - c$	$-c$	$-c$
×	0	0	0

(ii) Chase and ambush

Figure 1: Examples for  $m = 2$ ,  $b > c > 0$ . For each combination of *contribution* ( $+_1, +_2$ ) and *non-contribution* ( $\times$ ), entries give fitnesses for the row player.

we label the status quo strategy  $\times \in S_i$ . For now, we assume that actions other than  $\times$  exert (weakly) positive externalities (relative to  $\times$ ) on other individuals. This gives a public goods aspect to the game: a contribution of any form by  $i$  is at least as good for  $j$  as is non-contribution by  $i$ .

**(PG)** For all  $i, j \in M$ ,  $i \neq j$ ,  $s \in S$ , we have  $\pi_j(s_i, s_{-i}) \geq \pi_j(\times, s_{-i})$ .

For appropriately chosen status quo strategy profiles, condition (PG) is satisfied by all of the well known games in Table 1. Figure 1 gives two further examples that satisfy Condition (PG). In these examples, group size is two ( $m = 2$ ), there are two actions other than the status quo, and both group members are required to switch from the status quo in order to to gain some net benefit. Figure 1(i) is a stag hunt with two stags, and the hunters must both pursue the same stag in order to be successful. Figure 1(ii) represents a situation where there are two roles required for a successful hunt, such as when one hunter pursues the quarry and a second hunter lies in wait, ready to ambush the quarry when it flees from the first hunter.

From the status quo, the assumption that  $\underline{\times}$  is a Nash equilibrium implies that no individual acting alone can improve his payoff.<sup>6</sup> However, there exist

---

<sup>6</sup>We assume myopia. That is, individuals do not think, or rather act, beyond the implications of a direct adjustment to their strategy or the strategies of those with whom they share intentions. We are modeling early man, not bands of game theorists roving across the savannah. Note that even myopic payoff improvers are considerably more sophisticated than



Game	Description	Status quo $\times$
Threshold public goods	Players can <i>contribute</i> or <i>not contribute</i> . Contributing has a cost of $c$ . If at least $n$ players <i>contribute</i> , then every player receives a benefit $b > c$ .	The NE at which no players <i>contribute</i> .
Prisoner's dilemma	Players can <i>cooperate</i> or <i>defect</i> . Cooperating has a cost of $c$ for the cooperating player and generates a benefit of $b$ for every player. $b < c$ , $mb > c$ .	The unique NE, at which all players <i>defect</i> .
Trust game	Two players. Player 1 chooses an amount of some endowment to pass to Player 2. Any amount passed is tripled in value, following which Player 2 chooses some amount to pass back to Player 1.	The unique NE, at which Player 1 does not give anything to Player 2, and Player 2 never gives anything back to Player 1.
Centipede game	Two players. Firstly, Player 1 can choose <i>down</i> or <i>across</i> . If he chooses down, payoffs for Players 1,2 respectively are (2, 2). If he chooses across, then Player 2 can choose <i>down</i> or <i>across</i> . If he chooses down, payoffs are (1, 4). If he chooses <i>across</i> , payoffs are (3, 3).	The unique NE, at which Player 1 chooses <i>down</i> and Player 2 chooses <i>down</i> when given the opportunity.
Nash demand game	Each player demands a share of some surplus. If the sum of demands is less than or equal to the total surplus, then every player receives his demand. Otherwise all players receive nothing.	The NE at which every player demands the entire surplus.
Bertrand oligopoly	There is fixed demand $D$ for a good. The marginal cost of producing the good is $c$ . Each player (firm) $i$ chooses a price $p_i \geq c$ . If there are $r$ players charging the lowest price of all the prices, each of these players receives payoff $(p_i - c)D/r$ . Other players receive zero payoff.	The NE at which every player chooses $p_i = c$ .
Minimum effort game	Each player chooses an effort $e_i \geq 0$ . His payoff is then given by $b \min_{j \in M} e_j - ce_i$ , where $b > c > 0$ .	The NE at which every player chooses $e_i = 0$ .

Table 1: Games and status quo strategy profiles that satisfy (PG). Games are defined for an arbitrary number of players ( $m \geq 2$ ) unless explicitly stated otherwise.

opportunities for coalitions of players who can share their intentions to collaborate and adjust their actions together in order to obtain higher payoffs. Let the set of collaborative opportunities for a set  $T \subseteq M$  be

$$\mathcal{C}(T) = \left\{ s \in S : \begin{array}{l} \text{For all } i \in T, s_i \neq \times \text{ and } \pi_i(s) > \pi_i(\underline{x}). \\ \text{For all } i \notin T, s_i = \times. \end{array} \right\}.$$

That is,  $\mathcal{C}(T)$  gives the ways in which individuals in  $T$  can collaboratively adjust their strategies so that their own payoffs improve, leaving the strategies of individuals outside of  $T$  fixed. We assume that the game affords at least some prospect of collaboration. That is,  $\mathcal{C}(T) \neq \emptyset$  for at least some  $T \subseteq M$ . This is equivalent to  $\underline{x}$  not being a Strong Equilibrium in the sense of [Aumann \(1959\)](#). For the example in Figure 1(i),  $\mathcal{C}(M) = \{(+_1, +_1), (+_2, +_2)\}$  and for the example in Figure 1(ii),  $\mathcal{C}(M) = \{(+_1, +_2), (+_2, +_1)\}$ . The size of the smallest coalition that can benefit from collaboration is

$$n = \min_{T: \mathcal{C}(T) \neq \emptyset} |T|$$

Note that our assumption that  $\underline{x}$  is a Nash equilibrium implies that  $n \geq 2$  and the existence of at least some collaborative opportunity implies that  $n \leq m$ . For the games in Figure 1, both players must adjust their strategy from the status quo in order for them to gain, so  $n = 2$ .

We shall allow for the possibility that the behavior of individuals outside of a set  $T$  will alter as a consequence of  $T$  exploiting a collaborative opportunity. With this in mind, define the set of outcomes that could occur following a set of individuals  $T$  exploiting a collaborative opportunity.

$$\mathcal{C}^*(T) = \left\{ s^* \in S : \text{For some } s \in \mathcal{C}(T), s_i^* = s_i \text{ for all } i \in T. \right\}.$$

That is, were a set of individuals  $T$  to adjust their strategies according to some collaborative opportunity in  $\mathcal{C}(T)$ , following which the remainder  $M \setminus T$  of the individuals were to adjust their strategies in some way, then  $\mathcal{C}^*(T)$  is the set of strategy profiles that could be reached.

---

types of players - *cooperators*, *defectors* and so on, who only play a specific action.

## 2.2. Types and behavior

There are two types of individual, those who can share intentions and those who cannot. Those who can share intentions can collaboratively optimize when choosing their action. We refer to such individuals as SI types. Those who lack the cognitive ability to engage in such joint optimization we refer to as N types.

From a status quo at which everybody is playing  $\times$ , some set of  $n$  individuals within a group could gain by adjusting their actions. However, doing this might not lead to a Nash equilibrium. In fact, it may even be that playing an action other than  $\times$  is never individually rational, such as in an  $m$ -person Prisoner's Dilemma. However, if nobody is playing anything other than  $\times$ , then for any set of  $n$  individuals to remain playing  $\times$  is not collectively rational. We defer consideration of individuals who are sometimes in the mood for collaboration and are sometimes not until later in the paper. For now we assume that SI types are always willing to participate in collaborative decisions when opportunities present themselves.<sup>7</sup>

Fix an  $m$ -player game  $\Gamma$  as described in Section 2.1. Consider a group  $M$  of  $m$  individuals who encounter this problem. Let  $M_{SI} \subseteq M$  denote the set of SI type individuals within the group. Then we assume the outcome of the game will be given by some strategy profile  $s^*$  satisfying

- (C) (i) If for all  $T \subseteq M_{SI}$ ,  $\mathcal{C}(T) = \emptyset$ , then  $s^* = \underline{\times}$ .
- (ii) If there exists  $T \subseteq M_{SI}$  such that  $\mathcal{C}(T) \neq \emptyset$ , then select some set of SI type individuals  $T \subseteq M_{SI}$ ,  $\mathcal{C}(T) \neq \emptyset$ , according to some probability measure  $F_{M_{SI}, \Gamma}(\cdot)$ . Then let  $s^* \in \mathcal{C}^*(T)$  be chosen according to some probability measure  $G_{M_{SI}, T, \Gamma}(\cdot)$ .

That is, when there is no subset of SI types within the group that can exploit a collaborative opportunity, the outcome of the game is the status quo  $\underline{\times}$ . When there exists at least one subset of SI type individuals who can exploit at least

---

<sup>7</sup>That is to say, N types would play a Prisoner's Dilemma, whereas SI types would play a Prisoners' Dilemma, the difference being in the positioning of the respective apostrophes.

one collaborative opportunity, then some set of SI types will exploit some collaborative opportunity. Note that without specifying  $F_{M_{SI},\Gamma}(\cdot)$  and  $G_{M_{SI},T,\Gamma}(\cdot)$ , condition (C) is not a complete description of behavior. In particular, when there are multiple collaborative opportunities any of them could be taken with any probability, and the behavior of the players who are not involved in exploiting the collaborative opportunity is similarly arbitrary.

Some results will require a weak additional assumption on the behavior of the non-collaborating players. In words, this assumption says that any individual's response to collaboration can be construed as part of a better response (R) by *some* group of players outside of  $T$  to *either* (i) the status quo, *or* (ii) the strategy profile following the exploitation of a collaborative opportunity by  $T$ .

**(R)** If  $s^* \in \text{supp } G_{M_{SI},T,\Gamma}(\cdot)$ ,  $i \notin T$ ,  $s_i^* \neq \times$ , then letting  $s$  be the unique  $s \in \mathcal{C}(T)$  corresponding to  $s^* \in \mathcal{C}^*(T)$ , there exists  $R \subseteq M \setminus T$  such that  $i \in R$  and for all  $j \in R$ , either  $\pi_j(s_R^*, \underline{x}_{-R}) > \pi_j(\underline{x})$  or  $\pi_j(s_R^*, s_{-R}) > \pi_j(s)$ .

This condition is satisfied if responses to collaboration are individual best responses to collaboration, but is clearly much weaker than that. It allows, in the spirit of the topic under consideration, that collaborative behavior may in turn generate collaborative responses. It is intended to rule out responses to collaboration whereby individuals adjust their strategies to reduce their own payoffs.

### 2.3. Matching

We consider a population comprising unit mass of individuals, each of whom may be of SI or N type. Let the share of SI types in the population be  $x_{SI}$  and the share of N types be  $x_N$ . Let the population state be  $x = (x_{SI}, x_N)$ .

Each member of the population is matched to play  $\Gamma$  in a group of  $m$  individuals. We assume that the allocation of individuals to player positions in the game  $\Gamma$  is independent of type.<sup>8</sup> Given a population state  $x$ , and an indi-

---

<sup>8</sup>Otherwise it would be possible to give either type an advantage by giving them preferred access to player positions that correspond to higher payoffs.

vidual, let  $Z$  be a random variable denoting the number of SI types amongst the other  $m - 1$  individuals with whom the given individual is matched. We allow correlation between  $Z$  and the type of the individual concerned. Write  $Pr_x[Z = k | SI]$  and  $Pr_x[Z = k | N]$  as the probabilities that there are  $k$  SI type individuals amongst the other members of the group, conditional on a given individual being SI type and N type respectively. A *matching protocol* specifies these values for all  $x$  and all values of  $k$  from 0 to  $m - 1$ . We assume that  $Pr_x[Z = k | SI]$  and  $Pr_x[Z = k | N]$  are continuous in  $x_{SI}$ , and strictly positive for  $x_{SI}, x_N > 0$ .

Given this notation, any SI type has a probability  $Pr_x[Z = k - 1 | SI]$  of being in a group that includes exactly  $k$  SI types, including himself. Therefore, the mass of SI types in such groups equals  $x_{SI}Pr_x[Z = k - 1 | SI]$ . Any N type has a probability  $Pr_x[Z = k | N]$  of being in a group that includes exactly  $k$  SI types. Therefore, the mass of N types in such groups equals  $x_NPr_x[Z = k | N]$ . Now, any group with  $k$  SI types has  $m - k$  N types, so the ratio of SI type individuals in such groups to N type individuals in such groups must equal  $k/(m - k)$ . Noting that  $x_N = 1 - x_{SI}$ , we have the balance condition (B).<sup>9</sup>

$$(B) \quad \frac{x_{SI}Pr_x[Z=k-1|SI]}{(1-x_{SI})Pr_x[Z=k|N]} = \frac{k}{m-k} \text{ for } k = 1, \dots, m - 1.$$

An implication of (B) is that as  $x_{SI}$  approaches zero, the ratio of  $Pr_x[Z = k - 1 | SI]$  to  $Pr_x[Z = k | N]$  approaches infinity. Consequently, SI types find themselves in groups in which collaboration occurs infinitely more often than N types find themselves in such groups. Furthermore,  $Pr_x[Z = k | N]$  must approach zero for  $k \geq 1$  and  $Pr_x[Z = 0 | N]$  must approach unity.

---

<sup>9</sup>Given the preceding, one might ask why it is that this is stated as an condition, rather than merely as a consequence of any matching protocol. The reason for this is that although (B) is a logical consequence of matching, it is not a mathematical consequence. Consider groups of size two, with SI types making up a share of a quarter of the population. It is mathematically possible to match this quarter of the population one-to-one to the remaining three quarters of the population via a bijection. This would clearly be against the spirit of the model.

2.4. *The main theorem*

Given a game  $\Gamma$ , some behavioral rule satisfying (C), and some matching protocol, let  $f_{SI}(x)$ ,  $f_N(x)$  denote the fitnesses, that is the expected payoffs, of SI and N types respectively at population state  $x$ . For  $|M_{SI}| = k$ , denote the expected payoff (taken before player positions within the game are allocated) of SI types within the group by  $\pi_{SI}^k$ , and denote the expected payoff of N types within the group by  $\pi_N^k$ . Assume that  $F_{M_{SI},\Gamma}(\cdot)$  and  $G_{M_{SI},T,\Gamma}(\cdot)$  are such that these expectations exist. Then

$$f_{SI}(x) = \sum_{k=0}^{m-1} Pr_x[Z = k | SI] \pi_{SI}^{k+1}, \quad f_N(x) = \sum_{k=0}^{m-1} Pr_x[Z = k | N] \pi_N^k.$$

Note that  $f_{SI}(x)$  and  $f_N(x)$  depend continuously on the probabilities  $Pr_x[Z = k | SI]$  and  $Pr_x[Z = k | N]$ , which in turn are continuous in  $x_{SI}$ . Therefore,  $f_{SI}(x)$  and  $f_N(x)$  are continuous in  $x_{SI}$ . For some subpopulation with shares of SI and N types given by  $\tilde{x}$ , let  $f_{\tilde{x}}(x)$  be the average fitness of members of this subpopulation when the population state is  $x$ . That is,

$$f_{\tilde{x}}(x) := \tilde{x}_{SI} f_{SI}(x) + \tilde{x}_N f_N(x).$$

We use the concept of an evolutionarily stable state ([Taylor and Jonker, 1978](#)). An evolutionarily stable state is a state such that following the invasion of the population by a small population share  $\varepsilon$  of mutants, the non-mutant share of the population outperforms the invading mutants.

**Definition 2.1.** A state  $x^*$  is an evolutionarily stable state (ESS) if for any other state  $\tilde{x}$ , defining  $x_\varepsilon = (1 - \varepsilon)x^* + \varepsilon\tilde{x}$ , there exists  $\tilde{\varepsilon}$  such that

$$\text{For all } \varepsilon < \tilde{\varepsilon}, \quad f_{x^*}(x_\varepsilon) > f_{\tilde{x}}(x_\varepsilon).$$

An interior state  $x^*$  is an ESS if and only if  $f_{SI}(\cdot) - f_N(\cdot)$  is strictly decreasing at  $x^*$  and equal to 0. The extremal state  $x_{SI}^* = 0$  ( $x_{SI}^* = 1$ ) is an ESS if and only if  $f_{SI}(\cdot) - f_N(\cdot)$  is strictly negative (positive) in some open interval bounded below (above) by  $x_{SI}^*$ . This implies that, unless there exists some open interval

of  $x_{SI}$  on which  $f_{SI}(\cdot) - f_N(\cdot) = 0$ , at least one ESS must exist. Such examples can be constructed, but will necessarily be special.<sup>10</sup>

We are now in a position to state our main theorem. For any problem satisfying (PG) and any matching protocol satisfying (B), SI types will make up a positive share of the population at any evolutionarily stable state. Even when conditions are highly adverse to the evolution of shared intentions, SI types will still persist as a small share of the population, ready to expand and take a greater share as soon as conditions become more favorable.

**Theorem 1.** *If (C),(R),(B),(PG) hold, then  $x_{SI} > 0$  in any ESS.*

The reasoning behind the Theorem is as follows. (i) Firstly, collaboration (C) is a mutualistic act, not an altruistic act, therefore when SI types collaborate they improve their payoffs relative to the status quo, holding fixed the strategies of the non-collaborators. (ii) Secondly, the public goods condition (PG) ensures that any response by other players to the collaboration can only (weakly) increase the payoffs of the collaborators. (iii) Thirdly, (PG) ensures that collaboration also (weakly) increases the payoffs of non-collaborators. (iv) Fourthly, (R) together with (PG) ensures that the response to collaboration by non-collaborators (weakly) increases the payoffs of non-collaborators. (v) Finally, the balance condition (B) implies that when SI types are a small share of the population, any given SI type will find himself in a group in which collaboration occurs much more frequently than any given N type finds himself in such a group. Thus an SI type will enjoy the benefits of collaboration, as either a collaborator or a free rider, far more often than an N type will get the opportunity to free ride on the collaboration of others.

Note that when  $n = m$ , any collaboration that takes place will always involve every group member, so we no longer need to consider the response of non-collaborators to collaboration (steps (ii) and (iv) above), or the effect of

---

<sup>10</sup>To make specific statements about genericity requires consideration not only of  $\Gamma$ , but also of different behavioral rules satisfying (C), and different matching processes. This is sufficiently involved that it is omitted here.

collaboration on non-collaborators (step (iii) above). As these steps are the only steps in the proof of Theorem 1 that use conditions (PG) and (R), these conditions become unnecessary. Furthermore, as N types will never be members of a group in which collaboration occurs, step (v) and hence condition (B) becomes unnecessary. Therefore, SI types outperform N types and a result stronger than Theorem 1 holds, regardless of whether (R),(B) or (PG) hold.

**Corollary 1.** *If (C) holds and  $n = m$ , then a unique ESS  $x^*$  exists and  $x_{SI}^* = 1$ .*

One important class of games that does not always satisfy (PG) but is covered by Corollary 1 is the class of coordination games for which all symmetric pure strategy profiles are Nash equilibria and  $\underline{x}$  is Pareto dominated by one of the other equilibria. Another important class is the class of two player games with a Pareto dominated Nash equilibrium. This latter class includes the Cournot duopoly game, which does not satisfy (PG), and the two player Prisoner's Dilemma, which does satisfy (PG).

Now, consider a finitely or infinitely repeated game  $\bar{\Gamma}$  for which the stage game satisfies (PG) and has a nonempty set of collaborative opportunities. Let the status quo strategy profile  $\underline{x}$  of  $\bar{\Gamma}$  be the status quo strategies of  $\Gamma$  repeated each period, irrespective of history. Then, considering situations in which there is no response to collaboration by non-collaborators, we can use steps (i),(iii) and (v) of the proof of Theorem 1 outlined above to show the following result.

**Corollary 2.** *If  $\bar{\Gamma}$  is a repeated game with a stage game  $\Gamma$  satisfying (PG), and the status quo profile  $\underline{x}$  for  $\bar{\Gamma}$  is the history independent repetition of the status quo strategies of  $\Gamma$ , then if (C),(B) hold, and  $G_{M_{SI},T,\Gamma}(\mathcal{C}(T)) = 1$  whenever it is defined, then  $x_{SI}^* > 0$  in any ESS.*

Of course, if *any* game satisfies (PG) then Theorem 1 applies. Therefore, the additional contribution of Corollary 2 is limited to cases in which a repeated game  $\bar{\Gamma}$  does not satisfy (PG), but its stage game  $\Gamma$  does satisfy (PG).

Finally, as the state  $x$  is unidimensional, evolutionarily stable states correspond to strongly uninvadable states in the sense of [Bomze \(1991\)](#) and are thus



	+ <sub>1</sub>	×	
+ <sub>1</sub>	b - c	b - c	
×	b	0	
	+ <sub>1</sub>	×	

	+ <sub>1</sub>	×	
+ <sub>1</sub>	b - c	-c	
×	0	0	
	+ <sub>1</sub>	×	

Figure 2: Three player threshold public goods game. Any two players must contribute for the good to be provided.  $m = 3$ ,  $n = 2$ ,  $b > c > 0$ . For each combination of *contribution* (+<sub>1</sub>) and *non-contribution* (×), entries give fitnesses for the row player.

locally asymptotically stable under the replicator dynamic (Bomze and Weibull, 1995).

### 2.5. Example: threshold public goods problems

Let  $\Gamma$  be a threshold public goods game with threshold  $n \leq m$ . There are two strategies, contribute (+<sub>1</sub>) and don't contribute (×). If at least  $n$  individuals contribute then the good is provided, otherwise it is not provided. When the good is provided, every individual in the group obtains a benefit of  $b$ . When an individual contributes, he incurs a cost of  $c$ . Assume  $b > c > 0$ . The case of  $m = 3$ ,  $n = 2$  is shown in Figure 2. Let matching be non-assortative, such that a given individual's type does not affect the type distribution of the remaining  $m - 1$  individuals with whom he is matched. That is,

**(NA)**  $Pr_x[Z = k] := Pr_x[Z = k | SI] = Pr_x[Z = k | N]$  for all  $k, x$ .

Note that when combined with (B), (NA) implies binomial matching (Bin).

**(Bin)**  $Pr_x[Z = k] = \binom{m-1}{k} x_{SI}^k (1 - x_{SI})^{m-1-k}$ .

Assume that when  $|M_{SI}| \geq n$ , the realized strategy profile,  $s^*$ , always has exactly  $n$  SI types contributing.

**(EX)** If  $|M_{SI}| \geq n$ , then  $F_{M_{SI}, \Gamma}(\{T : |T| = n\}) = 1$  and  $G_{M_{SI}, T, \Gamma}(\mathcal{C}(T)) = 1$ .

Note that (EX) satisfies (R) and seems plausible, as there is no advantage to anyone from any additional individual contributing.<sup>11</sup> This implies that from the perspective of an SI type, if the good is provided and there are  $k$  other SI types in the group, he will contribute  $n/k + 1$  of the time. The average fitness of SI types in this setting is

$$f_{SI}(x) = \sum_{k=n-1}^{m-1} Pr_x[Z = k] \left( b - c \frac{n}{k+1} \right). \quad (1)$$

N types will benefit when the good is provided, but are never able to be part of a collaborative effort to provide the good. Therefore, the good will only be provided if at least  $n$  of the other  $m - 1$  players are SI types. When  $n = m$ , the fitness of an N type is always zero. When  $n < m$ , the fitness of an N type is

$$f_N(x) = Pr_x[Z \geq n] b. \quad (2)$$

When  $x_{SI} > 0$ , the fitness advantage (disadvantage when negative) of SI types over N types is

$$f_{SI}(x) - f_N(x) = Pr_x[Z = n - 1] \left( b - \sum_{k=n-1}^{m-1} \frac{Pr_x[Z = k]}{Pr_x[Z = n - 1]} c \frac{n}{k+1} \right) \quad (3)$$

This expression equals  $Pr_x[Z = n - 1](b - c) > 0$  when  $n = m$ . When  $n < m$ , then by showing that the term in brackets is decreasing in  $x_{SI}$ , positive for small positive values of  $x_{SI}$  and negative for values of  $x_{SI}$  close to 1, we prove the following proposition.

**Example 1.** *For threshold public goods games, when (C), (EX), (NA), (B) hold,*

---

<sup>11</sup>Note that our simple story of when the good is provided is borne out by more complex dynamic processes. If the number of SI types in a group is at least  $n$ , then under a coalitional better response dynamic with uniform mistakes (see, for example [Newton and Angus, 2015](#)), provision of the good is uniquely stochastically stable in the sense of ([Young, 1993](#)). If the number of SI types is strictly less than  $n - 1$ , then non-provision is uniquely stochastically stable. Note that a stochastic stability analysis of such problems under *individualistic* best response dynamics is given by [Myatt and Wallace \(2008\)](#), but results change when coalitional behavior is allowed.

- (i) *There is a unique ESS,  $x^*$ . If  $n < m$ , then  $x_{SI}^* \in (0, 1)$ , and if  $n = m$ , then  $x_{SI}^* = 1$ .*
- (ii) *From any mixed population such that  $x_{SI}, x_N > 0$ , the replicator dynamic converges to  $x^*$ .*
- (iii)  *$x_{SI}^*$  decreases in  $m$ , increases in  $n$ , increases in  $b/c$ .*
- (iv)  *$x_{SI}^* \rightarrow 0$  as  $b/c \rightarrow 1$  and  $x_{SI}^* \rightarrow 1$  as  $b/c \rightarrow \infty$ . In particular, note that  $x_{SI}^*$  may be greater or less than  $n/m$ .*

As threshold public goods games satisfy (PG), Theorem 1 applies and tells us that even when conditions are bad for collaboration, a minority of SI types will persist in the population. Example 1 shows that, for threshold public goods problems, such conditions are when  $m$  is large relative to  $n$  so that there are many free riders whenever the public good is provided, or when the benefit-cost ratio  $b/c$  is low. This minority of SI types can then expand when changes in the environment or technology lead to conditions which are more favourable for collaborative behavior. Such a change could be an increase in  $n$  caused by an increase in the availability of larger prey due to migration, or changes in the climate when moving between glacial and interglacial periods. Another example would be a reduction in  $c$  due to reduced risks from hunting due to improved technology providing better weapons.

### 2.6. Example: three player prisoner's dilemma

Consider the three player ( $m = 3$ ) prisoner's dilemma in Figure 3. If  $c \geq 2b$ , then all collaborative opportunities involve three players ( $n = 3$ ), so by Corollary 1 there is a unique ESS at  $x_{SI}^* = 1$ . At such an ESS, every individual in every matched group will always play the cooperative action  $+_1$ .

Now, consider a lower cost of cooperation. If  $c < 2b$ , then there exist collaborative opportunities for sets of two players ( $n = 2$ ). Now, as  $x_{SI} \rightarrow 1$ , by (B) we have that SI types will almost always be in groups containing three SI types. In such groups either all three players will collaborate and each will receive a

	$+_1$	$\times$	
$+_1$	$3b - c$	$2b - c$	
$\times$	$2b$	$b$	
	$+_1$	$\times$	

	$+_1$	$\times$	
$+_1$	$2b - c$	$b - c$	
$\times$	$b$	$0$	
	$+_1$	$\times$	

Figure 3: Three player prisoner's dilemma.  $m = 3$ ,  $3b > c > b$ . For each combination of *cooperate* ( $+_1$ ) and *defect* ( $\times$ ), entries give fitnesses for the row player.

payoff of  $3b - c$ , or any two of the three players will collaborate and the expected payoff of each individual will be  $2b - 2c/3$  as each is only a collaborator (hence paying the cost) with probability  $2/3$ . As  $3b - c > 2b - 2c/3$  we have that as  $x_{SI} \rightarrow 1$ ,  $f_{SI}(x)$  eventually drops below  $3b - c + \varepsilon$  for any  $\varepsilon > 0$ .

Now, under (NA),(B), as  $x_{SI} \rightarrow 1$  we have that any N type will almost always be matched with two SI types. These SI types will collaborate to cooperate and the N type will obtain a payoff of  $2b$ . That is,  $f_N(x) \rightarrow 2b$  as  $x_{SI} \rightarrow 1$ . As  $2b > 3b - c$ , we then have that for  $x_{SI}$  close to 1, N types obtain higher fitness than SI types. That is, there does not exist an ESS with  $x_{SI}^* = 1$ . At any ESS, there will always be N type individuals who play  $\times$  when they are matched to play the game.

So we see that reducing the cost of cooperation can lead to less collaboration and hence to less cooperation. There is a clear distinction between collaboration and cooperation, a distinction that will become even more clear in the next section where we shall drop condition (PG) and consider the possibility of collaborative behavior that has a negative effect on non-collaborators.

### 3. Collaboration with negative externalities

A plausible sounding conjecture would be that if collaborating players gain fitness from their collaboration following any response by the other players, then SI will evolve. This conjecture is false. When externalities from collaboration

	+ <sub>1</sub>	×
+ <sub>1</sub>	-4	1
×	-3	0
	+ <sub>1</sub>	×

	+ <sub>1</sub>	×
+ <sub>1</sub>	1	-1
×	0	0
	×	×

Figure 4: Three player Hawk-Dove game. Any two players must attack the remaining player for an attack to be successful.  $m = 3$ ,  $n = 2$ . For each combination of *hawk* ( $+_1$ ) and *dove* ( $\times$ ), entries give fitnesses for the row player.

are negative and SI types are disproportionately likely to match with other SI types, then negative externalities caused by collaborating SI types on other SI types can outweigh the benefits of collaboration. To see this, first formalize the condition, Profitable Collaboration (PC).

**(PC)** If  $T \subseteq M_{SI}$ ,  $\mathcal{C}(T) \neq \emptyset$ ,  $s^* \in \text{supp } G_{M_{SI}, T, \Gamma}(\cdot)$ , then  $\pi_i(s^*) > \pi_i(\underline{\times})$  for all  $i \in T$ .

This is a weaker condition than (PG). Any setup that satisfies (PG) will satisfy (PC), but the converse is not true.

Consider the three player game in Figure 4. We call this a three player Hawk-Dove game as any two players can exploit the third (steal his food) and it is never worthwhile for the third player to resist this. Thus there are three asymmetric Nash equilibria (in fact, Strong Equilibria) in which two players exploit the remaining player. However, there is also another Nash equilibrium at which all players play  $\times$ . The only available collaborative opportunity is for two SI type players to exploit the third (eg.  $s = (+_1, +_1, \times)$ ). When a pair of players  $T$  take such a collaborative opportunity  $s \in \mathcal{C}(T)$ , the exploited player would lose payoff by adjusting his strategy, so (R) implies that  $G_{M_{SI}, T, \Gamma}(s) = 1$ . Note that (PC) is satisfied. Now, the fitness of an N type is

$$f_N(x) = \underbrace{Pr_x[Z = 2 | N]}_{\text{Prob. of N type being exploited by SI types}} (-3) \underbrace{\xrightarrow{x_{SI} \rightarrow 0}}_{\text{by (B)}} 0.$$

The fitness of an SI type is

$$f_{SI}(x) = Pr_x[Z = 2 | SI] \underbrace{\left( \frac{2}{3}(1) + \frac{1}{3}(-3) \right)}_{\substack{\text{When all three are SI,} \\ \text{2/3 chance of being} \\ \text{an exploiter}}} + Pr_x[Z = 1 | SI](1),$$

so if  $\lim_{x_{SI} \rightarrow 0} Pr_x[Z = 2 | SI] > 3 \lim_{x_{SI} \rightarrow 0} Pr_x[Z = 1 | SI]$ , then  $\lim_{x_{SI} \rightarrow 0} f_{SI}(x)$  is bounded above by a number strictly below zero. That is, positive assortative matching can cause SI types to have lower fitness than N types, even when the share of SI types in the population is small.

Now, consider the case where even for SI types, the probability of encountering other SI types decreases as the share of SI types in the population goes to zero. This is the rare encounters in the limit (REL) condition.

**(REL)** For  $k \geq 1$ ,  $\frac{Pr_x[Z=k-1 | SI]}{Pr_x[Z=k | SI]} \rightarrow \infty$  as  $x_{SI} \rightarrow 0$ .

Note that under (B), the no assortativity condition (NA) implies (REL), but (REL) does not imply (NA). It turns out that when (REL) holds and collaboration is always profitable, then SI types will always make up a strictly positive share of the population at any ESS.

**Theorem 2.** *If (C),(B),(REL),(PC) hold, then  $x_{SI} > 0$  in any ESS.*

The intuition behind the Theorem is simple. Note that the balance condition (B) implies that for  $k \geq n$ ,  $Pr_x[Z = n - 1 | SI] / Pr_x[Z = k | N]$  approaches infinity as  $x_{SI}$  approaches zero. Similarly, (REL) implies that for  $k \geq n$ ,  $Pr_x[Z = n - 1 | SI] / Pr_x[Z = k | SI]$  approaches infinity as  $x_{SI}$  approaches zero. That is, when collaboration occurs it will usually be when there are exactly  $n$  SI types in the group. (PC) implies that these collaborators gain fitness from their collaboration, and there are no other SI types in the group to be affected by any negative externalities. Hence, SI types outperform N types for small, positive values of  $x_{SI}$ .

For the Hawk-Dove example of Figure 4, under (REL) we have that  $f_N = Pr_x[Z = 2 | N](-3) \leq 0$ , and  $f_{SI} = Pr_x[Z = 2 | SI] \left( \frac{2}{3}(1) + \frac{1}{3}(-3) \right) + Pr_x[Z =$

$1 | SI](1)$ , which (REL) implies is positive for small values of  $x_{SI}$ . Then  $f_N < 0$  and  $f_{SI} > 0$  for small values of  $x_{SI}$ . Once again, SI types proliferate when rare.

Finally, note that an effect of (REL) is that, for small  $x_{SI}$ , the usual absence of non-collaborating SI types when collaboration occurs makes the effect of the response of the non-collaborating players on themselves (step (iv) in the proof steps of Theorem 1 above) irrelevant, so we can drop (R) from the conditions, although (R) may, as in the above example, be a component in the satisfaction of (PC). However, as games satisfying (PG) satisfy (PC) for any  $G_{\dots, \Gamma}(\cdot)$ , if we add (REL) then we can simply drop (R) from Theorem 1 to obtain the following.

**Theorem 3.** *If (C),(B),(REL),(PG) hold, then  $x_{SI} > 0$  in any ESS.*

#### 4. A continuum of types ordered by likelihood of collaboration

Consider a model where instead of two types, we have a continuum of types, specifically the unit interval. Each time he faces a problem, any given individual of type  $\sigma \in [0, 1]$  will be in a *collaborative mood* with probability  $\sigma$ , and in an *individualistic mood* with probability  $1 - \sigma$ . An individual in an individualistic mood will behave as an N type and an individual in a collaborative mood will behave as an SI type. Let the state,  $x$ , be a probability measure on the Borel sets  $\mathcal{B}([0, 1])$ . This approach, modeling the same individual as sometimes collaborative and sometimes not, is that suggested by [Bacharach \(2006\)](#) for dealing with potential conflicts between individual and collective rationality. Any individual, when facing a problem as part of a group, will sometimes be driven by individual considerations and sometimes by collective considerations.<sup>12</sup> Define  $\bar{\sigma}(x) := \int_{[0,1]} \sigma x(d\sigma)$  as the probability that a randomly drawn individual from

---

<sup>12</sup>[Bacharach \(2006\)](#) thinks of individuals as sometimes reasoning individually and sometimes engaging in ‘team reasoning’. Our assumptions relate to behavior and not to reasoning per se, but our model can, should the reader wish, be interpreted as a model of the evolution of team reasoning, specifically what [Bacharach](#) refers to as *restricted team reasoning*, where at any given point in time, not every individual can team reason but those that can, recognize one another as such.

a population at state  $x$  is in a collaborative mood. Let  $Z$  be the number of individuals with whom a given individual is matched who are in a collaborative mood. We can adapt binomial matching to this setting.

$$\text{(Bin-}\sigma) \Pr_x[Z = k] = \binom{m-1}{k} (\bar{\sigma}(x))^k (1 - \bar{\sigma}(x))^{m-1-k}.$$

An evolutionarily stable state will not typically exist. The reason for this is that under binomial matching, at any interior ESS, any individual in the population must have equal expected fitness when he collaborates and when he does not collaborate. But then, from any state  $x$  such that  $x(\{0\}) \neq 1$ ,  $x(\{1\}) \neq 1$ , a small mutant subpopulation with type shares  $\tilde{x} \neq x$  could emerge such that  $\bar{\sigma}(\tilde{x}) = \bar{\sigma}(x)$ . That is, some of the mutants have increased  $\sigma$  and some have decreased  $\sigma$ , but the average remains the same as before. Such mutant invasions do not alter the expected fitness of any individual in the population. In particular, the mutants still obtain the same average fitness as non-mutants, so  $x$  cannot be evolutionarily stable. Consequently, we use the weaker concept of Neutral Stability (Maynard Smith, 1982). Write  $g_\sigma(x)$  for the fitness of type  $\sigma$  at state  $x$ . Note that the average fitness of a subpopulation of types distributed according to  $\tilde{x}$  when the state is  $x$  is now

$$g_{\tilde{x}}(x) := \int_{[0,1]} g_\sigma(x) \tilde{x}(d\sigma).$$

A neutrally stable state is then a population state such that following the invasion of the population by a small population share  $\varepsilon$  of mutants, the invaders do not do better than the non-mutants.

**Definition 4.1.** A state  $\hat{x}$  is a neutrally stable state (NSS) if for any other state  $\tilde{x}$ , defining  $x_\varepsilon = (1 - \varepsilon)\hat{x} + \varepsilon\tilde{x}$ , there exists  $\tilde{\varepsilon}$  such that

$$\text{For all } \varepsilon < \tilde{\varepsilon}, \quad g_{\tilde{x}}(x_\varepsilon) \geq g_{\hat{x}}(x_\varepsilon).$$

Now, by definition of the behavior of type  $\sigma$ , and comparing (Bin) and (Bin- $\sigma$ ), we have

$$g_\sigma(\cdot) = \sigma f_{SI}(\bar{\sigma}(\cdot)) + (1 - \sigma) f_N(\bar{\sigma}(\cdot))$$



where  $f_{SI}, f_N$  denote fitnesses in the two type model under (Bin), slightly abusing notation to write  $x_{SI}$  rather than  $x$  as the argument of  $f_{SI}, f_N$ . This gives

$$g_x(\cdot) = f_N(\bar{\sigma}(\cdot)) + \bar{\sigma}(x)(f_{SI}(\bar{\sigma}(\cdot)) - f_N(\bar{\sigma}(\cdot))).$$

That is,  $g_x(\cdot)$ , and specifically  $g_x(x_\varepsilon)$ , is monotonic in  $\bar{\sigma}(x)$ . This implies we only need to check robustness of any conjectured NSS to invasions of extreme types  $\sigma = 0$  and  $\sigma = 1$ . These types correspond to N and SI types of the two type model. Now, for the two type model under (Bin),(PC), as (Bin) implies (B) and (REL), by Theorem 2 there exists an ESS with a positive share of SI types. Therefore, letting  $x^*$  be an ESS of the two type model, we have that  $\hat{x}$  such that  $\hat{x}(0) = x_N^*, \hat{x}(1) = x_{SI}^*$ , is an NSS of the continuum type model. Furthermore, as at an NSS under binomial matching, fitness from collaboration and non-collaboration must be the same, the only factor that affects the fitness of any given type is the distribution over how many of his fellow group members are in a collaborative mood. But under (Bin- $\sigma$ ), this distribution is completely determined by  $\bar{\sigma}(\cdot)$ . Therefore, if  $x'$  is vulnerable to an invasion by mutants, and  $\bar{\sigma}(x') = \bar{\sigma}(x'')$ , then  $x''$  must be vulnerable to the same mutant invasion. That is, the only factor that determines whether a state  $x$  is an NSS is the value of  $\bar{\sigma}(x)$ .

**Theorem 4.** *If (C),(Bin- $\sigma$ ),(PC) hold, then at least one NSS of the continuum model exists. Under these conditions, a state  $x$  of the continuum model is an NSS if and only if  $\bar{\sigma}(x) = x_{SI}^*$  for some ESS  $x^*$  of the two type model under (C),(Bin),(PC). This implies that a monomorphic NSS  $\hat{x}$  exists, with  $\hat{x}(\hat{\sigma}) = 1$  and  $\hat{\sigma} = x_{SI}^*$ .*

## 5. Discussion: Cooperation, magical thinking, altruism, commitment and conditional cooperation.

Here we compare the collaborative sharing of intentions to some other modes of behavior that have been considered in the literature. The crucial distinction is that collaboration involves coordination in *how* actions are chosen rather than

	+ <sub>1</sub>	×
+ <sub>1</sub>	$b - c$	$-c$
×	$0$	$0$

(i) Fitnesses

	+ <sub>1</sub>	×
+ <sub>1</sub>	$b - c$	$b - c$
×	$0$	$0$

(ii) Magical thinker

	+ <sub>1</sub>	×
+ <sub>1</sub>	$b - c$	$-c/2$
×	$-c/2$	$0$

(iii) Altruist

Figure 5: Two player threshold public goods game.  $m = n = 2$ ,  $b > c > 0$ . For each combination of *contribution* (+<sub>1</sub>) and *non-contribution* (×), entries give, for the row player, his (i) fitnesses and his preferences when he is a (ii) magical thinker and (iii) altruist.

in the chosen actions themselves. Naturally, collaboration will often lead to efficient coordination, but the two things are not the same, as we saw in the three player Hawk-Dove game, where an efficient symmetric Nash equilibrium is destroyed by collaboration.<sup>13</sup> This is the reason we avoid the use of the term “cooperation” in describing jointly intentional strategic choice, as practitioners have become accustomed to using the word “cooperation” to describe a state of efficient coordination rather than its attainment.

### 5.1. Cooperators

There has been much consideration in the academic literature of situations where one symmetric action profile Pareto dominates all other symmetric action profiles. The action corresponding to such a profile is then described as the “cooperative” action. Individuals who always play such an action are called *cooperators* and those who play an action corresponding to some inefficient Nash equilibrium are called *defectors*. In the absence of assortative matching, when there are few cooperators in the population, they will rarely match with one another and will be outperformed by defectors. That is, cooperators do not proliferate when rare and there exists an ESS in which they are absent from the

---

<sup>13</sup>For more on this point, see [Newton and Angus \(2015\)](#), where it is shown how coordinated action choice by small groups within a population can slow convergence to a globally efficient action profile, even when all players have perfectly common interests.

population. For the threshold public goods game of Section 2.5, this has been formally shown by Pacheco, Santos, Souza and Skyrms (2009).

### 5.2. *Magical thinkers*

Magical Thinkers erroneously attribute causal powers to their own decisions (Elster, 1979). Consider symmetric games and those magical thinkers who behave as if their fellow group members will always take the same action as they take.<sup>14</sup> This implies that they will always choose the action corresponding to the most efficient of all symmetric action profiles. For the threshold public goods game of Section 2.5, the fitness of any given individual in the  $m = n = 2$  case for each combination of contribution (+<sub>1</sub>) and non-contribution (×) by the individual and his fellow group member is given in Figure 5(i). However, the magical thinker will act as if his fitness is given by Figure 5(ii). From this we see that if  $b - c > 0$ , then magical thinkers will behave identically to cooperators, and if  $b - c < 0$ , then magical thinkers will behave as defectors. Unlike cooperators and defectors, magical thinkers are not automata, but the ordering that determines their choice of action (their preferences) differs from the ranking given by their fitnesses. This is not the case for SI types, whose preferences (which are given at the level of the individual) are unaffected by their SI-ness but who may, in collaboration with other SI types, choose action profiles from a richer set of options. Their preferences are the same, but the menu is larger.

Furthermore, magical thinking does not resolve coordination problems with even the simplest of asymmetries. Going back to the Chase and Ambush game of Figure 1(ii), it is clear that any mode of reasoning that leads one player to play +<sub>1</sub> will also lead the other player to play +<sub>1</sub> and the players will fail to coordinate. The best that can be hoped for is a mixed strategy equilibrium, which is fine, but inefficient when compared to what can be achieved by SI types. Now, forgetting asymmetric coordination and considering asymmetric

---

<sup>14</sup>These types are behaviorally equivalent to the ‘Kantian’ types of Alger and Weibull (2013).

payoffs, consider amending the Two Stags game of Figure 1(i) so that the row player attains a payoff of  $2b - c$  from coordination on the first stag, and the column player attains a payoff of  $2b - c$  from coordination on the second stag. It is clear that any individualistic reasoning process followed by both players, even if it takes into account the payoffs of the other player, will fail to attain efficient coordination. Magical thinking is thus insufficient to achieve the outcomes achieved by SI types in such games. This is related to the discussion of conditional cooperation in Section 5.5.

### 5.3. Commitment

In discussions about an earlier draft of this work, it has been proposed by a reader that the key feature of SI types is the ability to commit to a non-individually rational strategy. This cannot be the case in general, as for some games, for example coordination games and threshold public goods games, SI types can collaborate in such a way that the resulting strategy profile satisfies individual rationality. For other games, such as prisoner's dilemmas, there does indeed exist a conflict between collective rationality and individual rationality. Considering the two player prisoner's dilemma, N types will always stick with the status quo action and defect, whereas SI types, if playing against another SI type, will collaborate to cooperate. So, from the status quo, one type (N) plays a myopic individualistic best response, whereas the other type (SI) plays a coalitional Pareto improving response when this is possible. Thus both individual and collective rationality are represented. There is no type, say an FalseSI type, that pretends to be an SI type and agrees to mutualistic collaboration with SI types, but in fact defects. Such a type would, of course, under non-assortative matching, multiply in a population of SI types playing prisoner's dilemmas. However, this is a second order question to that posed in the current paper. In the same way that the concept of the truth must exist before the concept of a lie can make sense, the ability to make collaborative decisions must exist prior to the ability to cheat one's collaborative partners. In any case, in many, probably most, cases of human collaboration, collectively rational decisions do

	+ <sub>1</sub>	×
+ <sub>1</sub>	3	0
×	4	1

(i)

	+ <sub>1</sub>	×
+ <sub>1</sub>	3	2
×	2	1

(i-a)

	+ <sub>1</sub>	×
+ <sub>1</sub>	3	-4
×	4	1

(ii)

	+ <sub>1</sub>	×
+ <sub>1</sub>	3	0
×	0	1

(ii-a)

Figure 6: For each combination of *cooperate* (+<sub>1</sub>) and *defect* (×), entries give, for the row player in two prisoner’s dilemmas, his fitnesses [(i),(ii)] or his preferences when he is an altruist [(i-a),(ii-a) corresponding to (i),(ii) respectively].

not conflict with individual rationality.

#### 5.4. Altruists

Similarly to those of magical thinkers, the preferences of altruists differ from those of a fitness maximizing individual. Altruists will, when given the opportunity, sacrifice some amount of their own fitness in order to increase the fitness of others. This can sometimes solve coordination problems. Consider utilitarian altruists whose preferences correspond to maximizing the average fitness of those playing a game. For the prisoner’s dilemma in Figure 6(i) this gives preferences as in Figure 6(i-a). Given these preferences, an altruist will act as a cooperator and so, as discussed in Section 5.1, altruism will not proliferate when rare in the absence of assortative matching. However, altruism may not even solve the coordination problem to begin with, even for prisoner’s dilemmas. For the prisoner’s dilemma in Figure 6(ii), a utilitarian altruist will still face the coordination problem of Figure 6(ii-a). A similar comment applies to threshold public goods problems (Figure 5).

#### 5.5. Conditional cooperators

Conditional cooperators identify the type of those with whom they interact and condition their action choice on this information (Hamilton, 1964a,b). This has been called a *green-beard* effect (Dawkins, 1976), as individuals with some observable characteristic - a “green beard”, behave cooperatively towards other

	$+_1$	$\times$
$+_1$	3, 3	0, 2
$\times$	2, 0	2, 2

(i)

	$+_1$	$\times$
$+_1$	1, 1	0, 0
$\times$	0, 0	2, 2

(ii)

	$+_1$	$\times$
$+_1$	3, 1	0, 0
$\times$	2, 0	2, 2

(iii)

Figure 7: Three games which are strategically equivalent from an individual perspective. For each combination of  $+_1$  and  $\times$ , entries give fitnesses for the row and column players respectively.

individuals with this characteristic. When there is a unique “cooperative action” ( $S_i = \{\times, +_1\}$  for all  $i \in M$ , (PG) holds), conditional cooperators who play  $+_1$  if and only if there are least  $n$  conditional cooperators in the group are similar to SI types and will proliferate when rare. This is the case for the threshold public goods model of Section 2.5 and for prisoner’s dilemmas.

However, to expand the concept of conditional cooperation so as to be applicable to a large variety of games is a non-trivial task. All of the games in Figure 7 are strategically equivalent from an individual perspective: given expectations of the opponent’s strategy (including mixed strategies), optimal strategies are the same in each of the cases. However, we do not want the conditional cooperator who switches to  $+_1$  when he is matched to a conditional cooperator in game (i) to do the same thing in games (ii) or (iii). Hence a reasonable definition of what a conditional cooperator should do in different games has to depend on his own payoffs and the payoffs of his opponent. Specifically, to mimic the outcomes obtained by SI types, the collective rationality of strategy profiles must be considered.

Furthermore, it is not clear how conditional cooperation should work when there are multiple opportunities for mutualistic collaboration. In particular, if collaborative opportunities are asymmetric, such as in the Chase and Ambush game in Figure 1(ii), something more than merely conditioning on the other players being conditional cooperator types is required. One possibility would be for there to exist multiple types of conditional cooperator. For example, there

could exist individuals with multiple shades of green beard, with the individual who sports the lighter shade of beard playing  $+_1$  and the darker individual playing  $+_2$ . What this conditioning is of course doing, is to implement asymmetric coordination of strategic adjustment.

We have seen in the above two paragraphs that to extend the idea of conditional cooperation beyond specific games, we need to introduce both considerations of collective rationality and of coordinated strategic adjustment. But these are exactly the elements - collective rationality, coordinated strategic choice, optimization - that are incorporated into our SI types, who are endowed with a comprehensive, multipurpose faculty that can be used for all such problems.

## **6. Conclusion: a modest proposal**

It has been shown that for broad classes of games we can expect at least some degree of agency to act at a collective level as if motivated by shared intentions. More specifically, we can expect to observe behavior that accords with some degree of agency being exercised at a collective level. The paper is silent as to how this collective agency is created, which as noted in the introduction, could be via explicit communication, tacit understanding, or team reasoning. Such an approach is not unusual to economics, where concepts such as the ‘firm’ and the ‘household’ are frequently used. It is clear that when decisions at a firm or household level are discussed, some degree of collective agency must be present, although the nature of this collective agency is not usually made explicit.

However, game theory in economics is in a weaker position. The most commonly used solution concept, Nash equilibrium, is habitually used without any explicit justification. Moreover, many of the Nash equilibria that occur in the literature are not Strong Equilibria; they are not robust to coalitional deviation, or to use the language of the current paper, from a status quo of such a Nash equilibrium, there exists an opportunity for collaboration. Therefore, an implication of the current work is that when using the concept of Nash equilibrium, an economist should ask whether the equilibrium is a Strong Equilibrium, and

if it is not, should carefully consider the extent to which joint agency might be expected to manifest itself in the problem under consideration. For example, does the problem satisfy (PG), or would a collaborative move away from the Nash equilibrium in question be likely to satisfy (PC)? How large would the gains be for collaborators? What would the externalities of collaboration be? Moreover, this proposal is not just predicated on the work here, but also on rich empirical evidence that the ability to share intentions and pursue mutually beneficial goals together with others is a basic human trait that cannot be ignored by any field that purports to scientifically consider human action.

## Appendix A. Proofs

Denote the average status quo payoff, by

$$\underline{\pi} := \frac{1}{m} \sum_{i=1}^m \pi_i(\underline{x}), \quad i \in M.$$

Define

$$\pi_N^{max} := \max_{k \geq n} \pi_N^k, \quad \pi_{SI}^{min} := \min_{k \geq n} \pi_{SI}^k.$$

*Proof of Theorem 1.* The average fitness of an N type is bounded above by

$$f_N(x) \leq \underbrace{Pr_x[Z < n | N]}_{\substack{\text{Prob. too few} \\ \text{SI types for} \\ \text{collaboration}}} \underline{\pi} + \underbrace{Pr_x[Z \geq n | N]}_{\substack{\text{Prob. enough} \\ \text{SI types for} \\ \text{collaboration}}} \pi_N^{max}.$$

The average fitness of an SI type is bounded below by

$$f_{SI}(x) \geq \underbrace{Pr_x[Z < n - 1 | SI]}_{\substack{\text{Prob. too few} \\ \text{SI types for} \\ \text{collaboration}}} \underline{\pi} + \underbrace{Pr_x[Z \geq n - 1 | SI]}_{\substack{\text{Prob. enough} \\ \text{SI types for} \\ \text{collaboration}}} \pi_{SI}^{min}.$$

Subtracting,

$$\begin{aligned} f_{SI}(x) - f_N(x) &= (f_{SI}(x) - \underline{\pi}) - (f_N(x) - \underline{\pi}) \\ &\geq Pr_x[Z \geq n - 1 | SI](\pi_{SI}^{min} - \underline{\pi}) - Pr_x[Z \geq n | N](\pi_N^{max} - \underline{\pi}). \end{aligned} \tag{A.1}$$



Now, (B) implies that for small enough  $x_{SI}$ ,

$$Pr_x[Z \geq n - 1 | SI] > Pr_x[Z \geq n | N] \left( \frac{\pi_N^{max} - \underline{\pi}}{\pi_{SI}^{min} - \underline{\pi}} \right), \quad (\text{A.2})$$

By (C) and (PG), when an SI type in any position in the game is in the set of collaborators, he gets a payoff strictly greater than the status quo payoff for that position. By (R) and (PG), when an SI type in any position in the game is not in the set of collaborators, he gets a payoff at least as high as the status quo payoff for that position. So when  $k \geq n$ , SI types always do at least as well as the status quo payoff and sometimes strictly improve upon it. Therefore  $\pi_{SI}^{min} > \underline{\pi}$ .

So  $\pi_{SI}^{min} - \underline{\pi} > 0$ . Together with (A.2), this implies that the RHS of (A.1) is greater than zero for small enough  $x_{SI}$ . That is,  $x_{SI} = 0$  cannot be an ESS, so any ESS must have  $x_{SI} > 0$ .  $\square$

*Proof of Corollary 1.* When  $n = m$ , unless  $M_{SI} = M$ ,  $s^* = \underline{x}$ . Therefore, for any  $x$ , the fitness of an N type is

$$f_N(x) = \underline{\pi}.$$

Now,  $n = m$  implies  $\mathcal{C}(M) \neq \emptyset$  and  $\mathcal{C}(M) = \mathcal{C}^*(M)$ , so (C) implies that  $\pi_{SI}^m > \underline{\pi}$ . Therefore, for  $x_{SI} > 0$ , the fitness of an SI type is bounded below by

$$f_{SI}(x) \geq Pr_x[Z < m - 1 | SI] \underline{\pi} + Pr_x[Z = m - 1 | SI] \pi_{SI}^m > \underline{\pi}.$$

$\square$

*Proof of Corollary 2.* The proof is identical to the proof of Theorem 1, except when it comes to showing that  $\pi_{SI}^{min} > \underline{\pi}$ . By (C) and  $G_{M_{SI}, T, \Gamma}(\mathcal{C}(T)) = 1$ , when an SI type in any position in the game is in the set of collaborators, he gets a payoff strictly greater than the status quo payoff for that position. By (PG) and  $G_{M_{SI}, T, \Gamma}(\mathcal{C}(T)) = 1$ , when an SI type in any position in the game is not in the set of collaborators, he gets a payoff at least as high as the status quo payoff for that position. So when  $k \geq n$ , SI types always do at least as well as the status quo payoff and sometimes strictly improve upon it. Therefore  $\pi_{SI}^{min} > \underline{\pi}$ .  $\square$

*Proof of Example 1.* The term in brackets in (3), simplified and divided by  $c$  equals

$$\frac{b}{c} - \sum_{k=n-1}^{m-1} \frac{n!(m-n)!}{(k+1)!(m-1-k)!} \left( \frac{x_{SI}}{1-x_{SI}} \right)^{k-(n-1)}, \quad (\text{A.3})$$

which is clearly strictly decreasing in  $x_{SI}$  when  $n < m$ , approaches  $b/c - 1 > 0$  as  $x_{SI} \rightarrow 0$ , and diverges to  $-\infty$  as  $x_{SI} \rightarrow 1$ . Therefore (A.3) equals zero and (3) crosses zero at some unique  $x^*$ , is strictly positive for all  $x$  such that  $x_{SI} \in (0, x_{SI}^*)$ , and is strictly negative for all  $x$  such that  $x_{SI} \in (x_{SI}^*, 1)$ . Therefore  $x^*$  is the unique ESS and the replicator dynamic converges to  $x^*$  from all  $x$  such that  $x_{SI} \in (0, 1)$ . Now, (A.3) increases in  $b/c$ ,  $n$  and decreases in  $m$ , so the value of  $x_{SI}$  at which (A.3) equals zero must increase or decrease respectively.

Fixing  $x_{SI} < 1$  and letting  $b/c \rightarrow \infty$ , the expression in (A.3) diverges to positive infinity, so  $x_{SI}^* \rightarrow 1$  as  $c \rightarrow 0$ . Conversely, if we fix  $x_{SI} > 0$  and let  $c \rightarrow b$ , then using expression (1) in the main text, we have

$$\begin{aligned} f_{SI}(x) &\rightarrow \sum_{k=n}^{m-1} Pr_x[Z = k] b \left( 1 - \frac{n}{k+1} \right) \\ &\leq Pr_x[Z \geq n] b \left( 1 - \frac{n}{m} \right) < Pr_x[Z \geq n] b = f_N(x), \end{aligned}$$

so  $x_{SI}^* \rightarrow 0$  as  $c \rightarrow b$ . □

*Proof of Theorem 2.*

$$f_N(x) - \underline{\pi} = \sum_{k=n}^{m-1} Pr_x[Z = k | N] (\pi_N^k - \underline{\pi})$$

and

$$f_{SI}(x) - \underline{\pi} = \sum_{k=n-1}^{m-1} Pr_x[Z = k | SI] (\pi_{SI}^{k+1} - \underline{\pi}),$$

giving

$$\begin{aligned} f_{SI}(x) - f_N(x) &= Pr_x[Z = n-1 | SI] \left( (\pi_{SI}^n - \underline{\pi}) + \right. \\ &\quad \left. \sum_{k=n}^{m-1} \frac{Pr_x[Z = k | SI]}{Pr_x[Z = n-1 | SI]} (\pi_{SI}^{k+1} - \underline{\pi}) - \frac{Pr_x[Z = k | N]}{Pr_x[Z = n-1 | SI]} (\pi_N^k - \underline{\pi}) \right). \end{aligned} \quad (\text{A.4})$$

Now, when  $|M_{SI}| = n$ , all possible collaborative opportunities involve all individuals in  $M_{SI}$ , so (PC) implies that  $(\pi_{SI}^n - \underline{\pi}) > 0$ . Furthermore, as discussed in the main body of the paper, (B) implies that for  $k \geq n$ ,  $Pr_x[Z = k | N] / Pr_x[Z = n - 1 | SI] \rightarrow 0$  as  $x_{SI} \rightarrow 0$ , and (REL) implies that for  $k \geq n$ ,  $Pr_x[Z = k | SI] / Pr_x[Z = n - 1 | SI] \rightarrow 0$  as  $x_{SI} \rightarrow 0$ . This implies that for small enough  $x_{SI} > 0$ , the right hand side of (A.4) is strictly positive, so any ESS must have  $x_{SI} > 0$ .  $\square$

*Proof of Theorem 3.* Immediate from the discussion prior to the statement of the Theorem in the main body of the text.  $\square$

*Proof of Theorem 4.* By definition,  $\hat{x}$  is an NSS if and only if for all  $\tilde{x}$  there exists  $\tilde{\varepsilon}$  such that for all  $\varepsilon < \tilde{\varepsilon}$ ,  $x_\varepsilon = (1 - \varepsilon)\hat{x} + \varepsilon\tilde{x}$ , we have

$$g_{\hat{x}}(x_\varepsilon) - g_{\tilde{x}}(x_\varepsilon) = (f_{SI}(\bar{\sigma}(x_\varepsilon)) - f_N(\bar{\sigma}(x_\varepsilon))) (\bar{\sigma}(\hat{x}) - \bar{\sigma}(\tilde{x})) \geq 0. \quad (\text{A.5})$$

If  $\bar{\sigma}(\hat{x}) = \bar{\sigma}(\tilde{x})$ , we have  $g_{\hat{x}}(x_\varepsilon) - g_{\tilde{x}}(x_\varepsilon) = 0$ . So (A.5) holds if and only if

(a) when  $\bar{\sigma}(\hat{x}) > \bar{\sigma}(\tilde{x})$ , there exists  $\tilde{\varepsilon}$  such that for all  $\varepsilon < \tilde{\varepsilon}$ ,

$$f_{SI}(\bar{\sigma}(x_\varepsilon)) - f_N(\bar{\sigma}(x_\varepsilon)) \geq 0,$$

and (b) when  $\bar{\sigma}(\hat{x}) < \bar{\sigma}(\tilde{x})$ , there exists  $\tilde{\varepsilon}$  such that for all  $\varepsilon < \tilde{\varepsilon}$ ,

$$f_{SI}(\bar{\sigma}(x_\varepsilon)) - f_N(\bar{\sigma}(x_\varepsilon)) \leq 0.$$

That is, at an interior NSS,  $f_{SI}(\cdot) - f_N(\cdot)$  must be weakly decreasing and equal to zero at  $\bar{\sigma}(\hat{x})$ . If  $\bar{\sigma}(\hat{x}) = 0$  ( $\bar{\sigma}(\hat{x}) = 1$ ) is an NSS, then  $f_{SI}(\cdot) - f_N(\cdot)$  must be weakly negative (weakly positive) on some open interval bounded below (above) by  $\bar{\sigma}(\hat{x})$ . As (Bin) implies that no part of  $f_{SI}(\cdot) - f_N(\cdot)$  is linear, these conditions are equivalent to  $f_{SI}(\cdot) - f_N(\cdot)$  being strictly decreasing and equal to zero at  $\bar{\sigma}(\hat{x})$ , or if  $\bar{\sigma}(\hat{x}) = 0$  ( $\bar{\sigma}(\hat{x}) = 1$ ),  $f_{SI}(\cdot) - f_N(\cdot)$  being strictly negative (strictly positive) on some open interval bounded below (above) by  $\bar{\sigma}(\hat{x})$ . These are the same conditions on  $\bar{\sigma}(\hat{x})$  as those placed on  $x_{SI}^*$  for an ESS of the two type model.  $\square$

## References

- Alger, I., Weibull, J.W., 2013. Homo moralis–preference evolution under incomplete information and assortative matching. *Econometrica* 81, 2269–2302.
- Alvard, M., 2001. Mutualistic hunting, in: Stanford, C., Bunn, H. (Eds.), *The early human diet: The role of meat*. Oxford University Press, Oxford, pp. 261–278.
- Alvard, M.S., Nolin, D.A., 2002. Rousseau’s whale hunt? *Current Anthropology* 43, 533–559.
- Ambrus, A., 2009. Theories of coalitional rationality. *Journal of Economic Theory* 144, 676 – 695.
- Angus, S.D., Newton, J., 2015. Emergence of shared intentionality is coupled to the advance of cumulative culture. *PLoS Comput Biol* 11, e1004587.
- Aumann, R., 1959. Acceptable points in general cooperative n-person games, in: Tucker, A.W., Luce, R.D. (Eds.), *Contributions to the Theory of Games IV*. Princeton University Press, pp. 287–324.
- Bacharach, M., 1999. Interactive team reasoning: a contribution to the theory of co-operation. *Research in economics* 53, 117–147.
- Bacharach, M., 2006. *Beyond individual choice: teams and frames in game theory*. Princeton University Press.
- Bernheim, B.D., Peleg, B., Whinston, M.D., 1987. Coalition-proof nash equilibria i. concepts. *Journal of Economic Theory* 42, 1–12.
- Bomze, I.M., 1991. Cross entropy minimization in uninvadable states of complex populations. *Journal of Mathematical Biology* 30, 73–87.
- Bomze, I.M., Weibull, J.W., 1995. Does neutral stability imply Lyapunov stability? *Games and Economic Behavior* 11, 173–192.

- Bowles, S., 2006. Group competition, reproductive leveling, and the evolution of human altruism. *Science* 314, 1569–1572.
- Bratman, M.E., 1992. Shared cooperative activity. *The Philosophical Review* 101, 327–341.
- Butterfill, S., 2012. Joint action and development. *The Philosophical Quarterly* 62, 23–47.
- Call, J., 2009. Contrasting the social cognition of humans and nonhuman apes: The shared intentionality hypothesis. *Topics in Cognitive Science* 1, 368–379.
- Choi, J.K., Bowles, S., 2007. The coevolution of parochial altruism and war. *Science* 318, 636–640.
- Dawkins, R., 1976. *The selfish gene*. revised edn. 1989 Oxford .
- Dekel, E., Ely, J.C., Yilankaya, O., 2007. Evolution of preferences. *The Review of Economic Studies* 74, 685–704.
- Elster, J., 1979. *Ulysses and the Sirens: Studies in Rationality and Irrationality*. Cambridge University Press.
- Eshel, I., Cavalli-Sforza, L.L., 1982. Assortment of encounters and evolution of cooperativeness. *Proceedings of the National Academy of Sciences* 79, 1331–1335.
- Farrell, J., Maskin, E., 1989. Renegotiation in repeated games. *Games and economic behavior* 1, 327–360.
- Fisher, R.A., 1930. *The Genetical Theory of Natural Selection*, ISBN 0198504403, variorum ed.(2000). Oxford University Press, USA.
- Gavrilets, S., 2014. Collective action and the collaborative brain. *Journal of The Royal Society Interface* 12.
- Gilbert, M., 1990. Walking together: A paradigmatic social phenomenon. *Midwest Studies in Philosophy* 15, 1–14.

- Gillies, D.B., 1959. Solutions to general non-zero-sum games. *Contributions to the Theory of Games* 4, 47–85.
- Gold, N., Sugden, R., 2007. Collective intentions and team agency. *The Journal of Philosophy* 104, 109–137.
- Güth, W., Kliemt, H., 1998. The indirect evolutionary approach: Bridging the gap between rationality and adaptation. *Rationality and Society* 10, 377–399.
- Haldane, J.B.S., 1932. *The causes of evolution*. Princeton University Press, 1990 ed.
- Hamilton, W., 1964a. The genetical evolution of social behaviour. i. *Journal of Theoretical Biology* 7, 1 – 16.
- Hamilton, W., 1964b. The genetical evolution of social behaviour. {II}. *Journal of Theoretical Biology* 7, 17 – 52.
- Hamilton, W.D., 1963. The evolution of altruistic behavior. *American naturalist* 97, 354–356.
- Kant, I., 1786. *What does it mean to orient oneself in thinking?*
- Maynard Smith, J., 1982. *Evolution and the Theory of Games*. Cambridge university press.
- Moll, H., Tomasello, M., 2007. Cooperation and human cognition: the Vygotskian intelligence hypothesis. *Philosophical Transactions of the Royal Society B: Biological Sciences* 362, 639–648.
- Myatt, D.P., Wallace, C., 2008. When does one bad apple spoil the barrel? an evolutionary analysis of collective action. *The Review of Economic Studies* 75, 499–527.
- Newton, J., 2012. Coalitional stochastic stability. *Games and Economic Behavior* 75, 842–54.

- Newton, J., Angus, S.D., 2015. Coalitions, tipping points and the speed of evolution. *Journal of Economic Theory* 157, 172 – 187.
- Pacheco, J.M., Santos, F.C., Souza, M.O., Skyrms, B., 2009. Evolutionary dynamics of collective action in n-person stag hunt dilemmas. *Proceedings of the Royal Society of London B: Biological Sciences* 276, 315–321.
- Robson, A.J., 1996. The evolution of attitudes to risk: Lottery tickets and relative wealth. *Games and economic behavior* 14, 190–207.
- Samuelson, L., 2001. Introduction to the evolution of preferences. *Journal of Economic Theory* 97, 225–230.
- Searle, J., 1990. Collective intentions and actions, in: Cohen, P.R., Morgan, J., Pollack, M. (Eds.), *Intentions in communication*. MIT Press, pp. 401–15.
- Smith, E.A., 2003. Human cooperation: Perspectives from behavioral ecology, in: P.Hammerstein (Ed.), *Genetic and cultural evolution of cooperation*. MIT Press, pp. 401–427.
- Sosis, R., Feldstein, S., Hill, K., 1998. Bargaining theory and cooperative fishing participation on Ifaluk atoll. *Human Nature* 9, 163–203.
- Sugden, R., 2000. Team preferences. *Economics and Philosophy* 16, 175–204.
- Taylor, P.D., Jonker, L.B., 1978. Evolutionary stable strategies and game dynamics. *Mathematical biosciences* 40, 145–156.
- Tomasello, M., 2014. *A natural history of human thinking*. Harvard University Press.
- Tomasello, M., Carpenter, M., 2007. Shared intentionality. *Developmental science* 10, 121–125.
- Tomasello, M., Carpenter, M., Call, J., Behne, T., Moll, H., 2005. Understanding and sharing intentions: The origins of cultural cognition. *Behavioral and brain sciences* 28, 675–691.

- Tomasello, M., Herrmann, E., 2010. Ape and human cognition what's the difference? *Current Directions in Psychological Science* 19, 3–8.
- Tomasello, M., Rakoczy, H., 2003. What makes human cognition unique? from individual to shared to collective intentionality. *Mind & Language* 18, 121–147.
- Trivers, R.L., 1971. The evolution of reciprocal altruism. *Quarterly review of biology* 46, 35–57.
- Tuomela, R., Miller, K., 1988. We-intentions. *Philosophical Studies* 53, 367–389.
- Velleman, J.D., 1997. How to share an intention. *Philosophy and Phenomenological Research: A Quarterly Journal* 57, 29–50.
- Vygotsky, L.S., 1980. *Mind in society: The development of higher psychological processes*. Harvard University Press.
- Wilson, D.S., Dugatkin, L.A., 1997. Group selection and assortative interactions. *American Naturalist* 149, 336–351.
- Wobber, V., Herrmann, E., Hare, B., Wrangham, R., Tomasello, M., 2014. Differences in the early cognitive development of children and great apes. *Developmental Psychobiology* 56, 547–573.
- Young, H.P., 1993. The evolution of conventions. *Econometrica* 61, 57–84.