

# Team Reasoning and a Rank-Based Function of Team's Interests\*

Jurgis Karpus<sup>†</sup>      Mantas Radzvilas<sup>‡</sup>

April 2015

## Abstract

Orthodox game theory is sometimes criticized for its failure to single out intuitively compelling solutions in certain types of interpersonal interactions. The theory of team reasoning provides a resolution in some such cases by suggesting a shift in decision-makers' mode of reasoning from individualistic to reasoning as members of a team. The existing literature in this field discusses a number of properties for a formalized representation of team's interests to satisfy: Pareto efficiency, successful coordination of individuals' actions and the notion of mutual advantage among the members of a team. For an explicit function of team's goals a reference is sometimes made to the maximization of the average of individuals' personal payoffs, which meets the Pareto efficiency and (in many cases) coordination criteria, but at times fails with respect to the notion of mutual advantage. It also relies on making interpersonal comparisons of payoffs which goes beyond the standard assumptions of the expected utility theory that make numerical representations of individuals' preferences possible. In this paper we propose an alternative, rank-based function of team's interests that does not rely on interpersonal comparisons of payoffs, axiomatizes the notion of mutual advantage and satisfies the weak Pareto efficiency and (in many cases) coordination criteria. We discuss its predictions using a number of examples and suggest a few possibilities for further research in this field.

## 1 Introduction

The standard rational choice theory is sometimes criticized for its inability to single out what at times appears to be the only obvious choice to make in

---

\*This is a work-in-progress paper, a similar version of which was presented at the UECE Lisbon Meetings 2014: Game Theory and Applications conference held at Lisboa School of Economics and Management (ISEG) on November 6–8, 2014, as well as at the Warwick Economics PhD Conference 2015 held at the University of Warwick on February 23–24, 2015.

<sup>†</sup>King's College London, Department of Philosophy (e-mail: jurgis.karpus@kcl.ac.uk).

<sup>‡</sup>London School of Economics, Department of Philosophy, Logic and Scientific Method (e-mail: m.radzvilas@lse.ac.uk).

games with multiple Nash equilibria. An example is the simple Hi-Lo game, in which two players independently and simultaneously choose one from a pair of available options: *Hi* or *Lo*. If both choose *Hi*, they get a payoff of 2 each. If both choose *Lo*, they get a payoff of 1 each. If one chooses *Hi* while the other chooses *Lo*, they both get 0. The game is illustrated in Figure 1, where one of the players chooses between two options identified by rows and the other—by columns. The numbers in each cell represent payoffs to the row and the column players respectively.

	<i>Hi</i>	<i>Lo</i>
<i>Hi</i>	2, 2	0, 0
<i>Lo</i>	0, 0	1, 1

Figure 1: The Hi-Lo game

The standard theory predicts that rational players will choose strategies that together constitute a Nash equilibrium, in which each player’s strategy is a best response to the strategies chosen by all other players. Here this means for both players to play *Hi* or for both to play *Lo*: (*Hi*, *Hi*) and (*Lo*, *Lo*) are Nash equilibria<sup>1</sup>. Yet (*Lo*, *Lo*) does not intuitively strike as a rational outcome in this game. It is true that if one player expected the other to play *Lo*, then choosing *Lo* would be his or her best response to the other player’s choice. In other words, choosing *Lo* would be the rational thing to do. However, it would be odd if anyone formed an expectation that a rational individual would play *Lo* in the first place. Experimental results support this by revealing that over 90% of the time people do opt for *Hi* in this game.<sup>2</sup>

This prompted the emergence of the theory of team reasoning which suggests that certain features of the context in which interdependent decisions are made may trigger a shift in peoples’ mode of reasoning from individualistic best-response reasoning to reasoning as members of a team where a group of individuals act together in the attainment of some common goal.<sup>3</sup> By identifying this goal with the maximization of the average of decision-makers’ personal payoffs the theory can be operationalized to render *Hi* to be the only rational choice in the Hi-Lo game for anyone who reasons as a member of a team.

---

<sup>1</sup>These are Nash equilibria in pure strategies. There is a third equilibrium in mixed strategies, in which players randomize between the two available options with certain probabilities. Here and in the rest of this paper we focus solely on equilibria in pure strategies.

<sup>2</sup>See Bardsley et al. (2010) who, among a number of other games, report experimental results from two versions of the Hi-Lo game where the outcome (*Hi*, *Hi*) yields a payoff of 10 while the outcome (*Lo*, *Lo*) yields a payoff of 9 or 1 to both players.

<sup>3</sup>For early developments of this theory see Sugden (1993, 2000, 2003) and Bacharach (1999, 2006). For some of the more recent work see Gold and Sugden (2007a,b), Sugden (2011, 2015) and Gold (2012).

Similarly, this allows to explain cooperation in the widely discussed Prisoner’s Dilemma and the Stag Hunt games.

While we agree with the idea that peoples’ mode of reasoning may sometimes undergo a shift from individualistic to reasoning as members of a team, we criticize the identification of the team’s interests with the maximization of the average of individuals’ personal payoffs. We do this for two reasons. First, it relies on making interpersonal comparisons of individual players’ payoffs, which goes beyond the standard assumptions of the expected utility theory. Second, it may advocate a complete self-sacrifice of some individuals for the benefit of others or possibly one sole member of a team. In this paper we present an alternative function for representing team’s interests that does not rely on making interpersonal comparisons of payoffs and makes participation in team play conditional on it bringing about at least some benefit to every member of the team.

The rest of this paper is structured as follows. In section 2 we discuss the theory of team reasoning in more detail and present how it is sometimes operationalized to render *Hi* to be the only rational choice in the Hi-Lo game and to explain cooperation in the Prisoner’s Dilemma and the Stag Hunt games. We also explain that the notion of team reasoning cannot be represented by a transformation of individuals’ personal payoffs in games. In section 3 we discuss why the maximization of the average of individuals’ personal payoffs may not be a good representation of team’s interests. In section 4 we present an alternative rank-based function<sup>4</sup> and illustrate its predictions using a variety of examples. With section 5 we conclude and suggest a few possible directions for further research.

## 2 Team Reasoning

When a person reasons individualistically, he or she focuses on the question “what it is that *I* should do in order to best promote *my* interests?”. The answer to this question identifies a strategy that is associated with the highest expected personal payoff *to the individual*, given his or her beliefs about the actions of others. This is what is meant by individualistic best-response reasoning underlying the identification of Nash equilibria in games. When a person reasons as a member of a team, he or she focuses on the question “what it is that *we* should do in order to best promote *our* interests?”. The answer to this question identifies a set of strategies—one strategy for each of the interacting individuals<sup>5</sup>—that leads to the attainment of the best possible outcome *for the*

---

<sup>4</sup>By *function* we mean something similar to what the rational choice theory refers to as a *choice function* that takes the set of the available actions to an individual, the structure of a game and the individual’s beliefs about others’ behaviour as inputs and produces a set of (rational) actions as an output. A slight difference in the case discussed in this paper is that the output of the function is a set of (rational) outcomes rather than actions.

<sup>5</sup>Strictly speaking this need not necessarily be the case, since not all individuals in a given strategic interaction may be reasoning as members of a team. There are variants of the theory of team reasoning that consider such scenarios. For an overview see Gold and Sugden (2007a).

*group of individuals* acting together as a team. As explained by Gold and Sugden (2007a) ‘when an individual reasons as a member of a team, she considers which combination of actions by members of the team would best promote the team’s objective, and then performs her part of that combination’.

In the existing literature on the theory of team reasoning that attempts to propose an explicit function for representing the interests of a team a reference is sometimes made to the maximization of the average of individuals’ personal payoffs. This suggestion can be found in Bacharach (1999, 2006) and, in later theoretical developments, it was made by Smerilli (2012). In empirical studies it was adopted by Colman et al. (2008, 2014). It is not the case that everybody operationalizes the theory of team reasoning using this function. Sugden (1993, 2000, 2003) as well as Gold and Sugden (2007a,b) do not endorse any explicit function and Gold (2012) suggests that team’s interests do not have to be represented by the maximization of average payoffs. Crawford et al. (2008) and Bardsley et al. (2010) interpret team reasoning as a search for a decision rule that would resolve certain types of coordination problems in a mutually beneficial way, but do not propose a specific function as well. Sugden (2011, 2015) suggests that a function representing team’s interests should incorporate the notion of mutual advantage among the members of a team and proposes a way of identifying mutually advantageous possibilities in games. We adopt the latter idea and expand it by axiomatizing mutually advantageous play in our construction of the rank-based function of team’s interests in section 4 below. With the exception of Sugden’s proposal to incorporate the notion of mutual benefit in formal representations of team’s interests, however, the literature on the theory of team reasoning that is known to us has not yet presented a formal alternative to the maximization of average payoffs as an explicit representation of team’s goals.

One reason why the maximization of the average of individuals’ payoffs is attractive is that it ensures Pareto efficiency<sup>6</sup> of any outcome that is selected by a team. It is easy to see that the outcome that best promotes this objective in the Hi-Lo game is  $(Hi, Hi)$ . Thus, when an individual reasons as a member of a team, he or she identifies the outcome  $(Hi, Hi)$  as uniquely optimal for the team and individually chooses  $Hi$ —his or her part in the attainment of this outcome.

Consider now the Prisoner’s Dilemma game illustrated in Figure 2A, in which two players independently and simultaneously decide whether to cooperate (play  $C$ ) or defect (play  $D$ ). The game has a unique Nash equilibrium:  $(D, D)$ . This is because, irrespective of what the other player is going to do, it is always better to play  $D$  from an individual’s personal point of view. Individualistic reasoning, thus, leads to a socially suboptimal, Pareto inefficient outcome, since the outcome  $(D, D)$  yields lower payoffs to both players than does the outcome  $(C, C)$ . This is the main reason why this game is so widely discussed in social

---

<sup>6</sup>An outcome of a game is Pareto efficient if there exists no other outcome in which somebody’s payoff could be increased without making anyone else worse off. In the Hi-Lo game the outcome  $(Lo, Lo)$  is Pareto inefficient since both players are better off in the outcome  $(Hi, Hi)$ , which is the only Pareto efficient outcome in this game.

sciences.

	<i>C</i>	<i>D</i>	
<i>C</i>	2, 2	0, 3	
<i>D</i>	3, 0	1, 1	
	<i>A</i>		

	<i>S</i>	<i>H</i>
<i>S</i>	2, 2	0, 1
<i>H</i>	1, 0	1, 1
	<i>B</i>	

Figure 2: The Prisoner’s Dilemma (A) and the Stag Hunt (B) games

Experimental results suggest that in a one-shot version of the Prisoner’s Dilemma game (i.e. when it is played once) people tend to cooperate about 50% of the time.<sup>7</sup> Notice that the outcome (*C*, *C*) uniquely maximizes the average of individuals’ personal payoffs. As such, the theory of team reasoning operationalized using the maximization of average payoffs suggests that in the Prisoner’s Dilemma game some people reason as members of a team while others reason individualistically.

In a similar way the theory can explain cooperation in the Stag Hunt game illustrated in Figure 2B. This game has two Nash equilibria: (*S*,*S*) and (*H*,*H*). However, hunting hare (playing *H*) guarantees a payoff of 1 irrespective of what the other player does, whereas the attainment of the high payoff from hunting stag (playing *S*) crucially depends on the cooperation of the other party. Here experimental results suggest that in a one-shot version of this game people tend to choose *S* slightly more than 60% of the time.<sup>8</sup> We will discuss this result in more detail in section 4.

## 2.1 Team Reasoning vs. Transformation of Personal Payoffs

An important point stressed by many game theorists is that payoff structures of games have to fully capture everything that is motivationally important in individuals’ evaluations of the possible outcomes of those games. For example, imagine that two people are playing the Prisoner’s Dilemma game in terms of monetary payoffs as illustrated in Figure 3A. Suppose that the row player is a pure altruist when it comes to decisions involving money—i.e. he or she always

<sup>7</sup>Cooperation tends to decrease with repetition, however (i.e. when people play the same game a number of times). See Ledyard (1995) for a survey of experimental results from *public goods* games, which involve more than two players but otherwise are very similar in their structure to the two-player Prisoner’s Dilemma.

<sup>8</sup>The proportion of people choosing *S* changes with repetition: in some cases it increases while in others it decreases. This seems to, at least partially, depend on the specific payoff structure of the played Stag Hunt game—the extent of risk involved in playing *S* and the extent of risklessness in playing *H*. See Battalio et al. (2001).

wants to maximize the other player’s monetary gain. Suppose also that the column player prefers to maximize his or her personal monetary payoff, but is extremely averse to inequitable distributions of gains among individuals. Assume that an outcome resulting in unequal monetary gains for the two players is just as good for him or her as gaining nothing. The correct representation of the true motivations of these individuals transforms the monetary Prisoner’s Dilemma game into the game illustrated in Figure 3B. This transformed game has a unique Nash equilibrium  $(C, C)$  in which both players rationally cooperate following individualistic best-response reasoning in the attainment of their personal goals.

	<i>C</i>	<i>D</i>		<i>C</i>	<i>D</i>
<i>C</i>	£2, £2	£0, £3	<i>C</i>	2, 2	3, 0
<i>D</i>	£3, £0	£1, £1	<i>D</i>	0, 0	1, 1
	<i>A</i>			<i>B</i>	

Figure 3: The Prisoner’s Dilemma game in terms of £ (A) and its transformation using payoffs that represent the true motivations of both players (B)

We agree with the idea that payoff structures of games have to accurately capture the true motivations of interacting individuals. It is important to note, however, that a possible shift in a decision-maker’s mode of reasoning from individualistic to reasoning as a member of a team cannot be captured by a transformation of that individual’s personal payoffs in the same way as it was done in the case of altruism and inequity-aversion above. To see this consider again the Hi-Lo game illustrated in Figure 1. Suppose that the row player reasons as a member of a team and adopts the maximization of the average of individuals’ personal payoffs as the team’s objective. Replacing his or her personal payoff numbers with averages of the two players’ payoffs in each outcome does not transform the original Hi-Lo game into anything different. This is because the personal payoffs of the two players already match the average of their personal payoffs in each outcome of this game to begin with. As a result, the transformed game would still have two Nash equilibria and two rational outcomes—exactly what the theory of team reasoning was developed to contest.

The difference between an individualistic mode of reasoning and reasoning as a member of a team lies not in how an individual personally values each outcome of a game, but in the way he or she reasons when choosing among the available actions. When a person reasons individualistically, he or she chooses an action based on his or her belief about what the other player is going to do and which outcome—and personal payoff—the chosen action would subsequently yield. When a person reasons as a member of a team, on the other hand,

he or she first identifies an outcome of the game that best fulfills the team's objective and then performs his or her part in the attainment of that outcome. In other words, somebody who reasons individualistically chooses an action that maximizes his or her expected personal payoff, whereas somebody who reasons as a member of a team chooses an action that is associated with an outcome that best fulfills the team's objective.

This underlines two key assumptions on which our approach is based. First, we assume that individuals' personal payoffs represent their true motivations in games. Second, we assume that the payoff structures of games are commonly known by all the interacting decision-makers. The latter point rules out the possibility that a mere shift in an individual's mode of reasoning—in addition to changing the way that individual reasons—changes the way he or she personally values each outcome in the considered games. If that were not the case, these interactions would cease to be games of complete information about each other's payoffs and would take us further away from orthodox game theory and the type of games we analyze in this paper.

Another important point to note is that reasoning as a member of a team does not imply and is not implied by the sharing of the attained payoffs among the members of a team. In other words, the attained payoffs are not transferable from one player to another. If players are able to (or are going to) share their personal gains, this has to be reflected in the payoff structures of the played games to begin with in order to correctly capture the players' true motivations. In this light, equal sharing of combined payoffs in the Hi-Lo game would leave the payoff structure of the original game unchanged. It would, however, change the payoff structure of any game that contained outcomes yielding unequal distributions of payoffs among players (as is the case, for example, in the Prisoner's Dilemma and the Stag Hunt games).

## **2.2 Team Reasoning and the Team's Interests: Two Separate Questions**

The theory of team reasoning, as introduced above, has to answer two separate but equally important questions. First, it needs to specify clear and testable circumstances under which individuals' mode of reasoning may undergo a shift from individualistic to reasoning as members of a team. Second, in cases when individuals do reason as members of a team, it needs to specify what they take the team's interests to be upon deciding on what courses of action to take. In this paper we predominantly focus on the latter question and turn to it next.

## **3 Team's Interests: Not the Average of Personal Payoffs**

As mentioned in the previous section, the literature on the theory of team reasoning that attempts to propose an explicit function for representing the interests of a team sometimes refers to the maximization of the average of individual

decision-makers' personal payoffs. We criticize this suggestion and believe that, if people do reason as members of a team in certain situations, they are unlikely to adopt the maximization of average payoffs as a guide when thinking about what it is that they should do. Our criticism is based on two points.

First, the maximization of the average of personal payoffs may, in certain situations, advocate a complete sacrifice of some individuals' personal interests for the benefit of others or possibly one sole member of a team—a consequence which we believe to be intuitively problematic. To see this, consider a slightly amended version of the Prisoner's Dilemma game illustrated in Figure 4A. The only difference from its original version is the slightly higher payoff to the row player from defection when the column player cooperates. In this game, the maximization of the average of players' personal payoffs would identify the team's objective with the attainment of the outcome  $(D,C)$ . It would thus prescribe a complete sacrifice of the column player's personal interests with the row player reaping all the benefits from team play.

	<i>C</i>	<i>D</i>		
<i>C</i>	2, 2	0, 3		
<i>D</i>	5, 0	1, 1		
			<i>A</i>	

	<i>L</i>	<i>R</i>
<i>U</i>	10, 1	0, 0
<i>D</i>	4, 4	1, 9
		<i>B</i>

Figure 4: The amended Prisoner's Dilemma (A) and the Chicken (B) games

For a slightly different example consider the game illustrated in Figure 4B, which is a particular version of the game known as the Chicken. This game has two Nash equilibria:  $(U,L)$  and  $(D,R)$ . Here the maximization of average payoffs suggests that team's interests would be best fulfilled with the attainment of the outcome  $(U,L)$ . In this case it does not prescribe a complete self-sacrifice to the column player—the outcome  $(U,L)$  is not the worst possible outcome for him or her in this game—but it does not advance the column player's personal interests anywhere far from just that.

Our suggestion is that, if some people's mode of reasoning does undergo a shift from individualistic to reasoning as members of a team in the amended Prisoner's Dilemma and the Chicken games above, the team's objective would be the attainment of the outcomes  $(C,C)$  and  $(D,L)$  respectively. More importantly, we suggest that no individual would willingly subscribe to team play if there was no personal gain for him or her from doing so, which precludes the possibility of a self-sacrifice. Our intuition is shared by Sugden (2015) who suggests that, in team play, 'each [player is] choosing his or her component of the joint action with the intention of achieving mutual benefit'.

This is not to say that self-sacrifice does not exist. It is evident that we



are often willing to sacrifice our personal material gains for the well-being of our loved ones. We suggest, however, that such motivational factors need to be fully captured by individuals' personal payoffs associated with the possible outcomes in games before the different modes of reasoning are considered. And, if team reasoning is a mode of reasoning that individual decision-makers may adopt in one-shot interactions with potentially complete strangers—cases that we consider in this paper—then self-sacrifice of individuals is unlikely.

The essence of this argument rests on an idea somewhat similar to the notion of the separateness of persons used to criticize utilitarianism. Utilitarianism, which focuses on the maximization of the aggregate well-being of a society while ignoring the distribution of the attained well-being among the individual members of that society, is said to be insufficiently sensitive to the separateness of individuals and the advancement of their personal well-being. In a similar fashion it can be argued that any aggregative function of team's interests that fails to take into account the distribution of payoffs among the members of a team is susceptible to being rejected for failing to respect the separateness of the interacting players. So long as there is space to make an objection on such grounds, any individual decision-maker could reasonably refuse taking part in the prescribed team play when their personal gains from doing so are insufficiently addressed.<sup>9</sup>

An important implication of the first point of our criticism is that we remain committed to the idea that all that matters for individual decision-makers are their personal motivations in the considered games. In this sense, the use of the term "team" in the interpretation of the theory of team reasoning in this paper is meant to be very loose. Most importantly, it is not meant to carry any psychological connotations that may be present in various other interpretations of team work and team members' duties (e.g. in sport, work groups, neighbourhoods or faculties) which may lead to individuals' abandonment of their personal aspirations (even when these aspirations take into account the well-being of others). In certain types of interpersonal interactions, however, the decision-makers' personal interests may be advanced further relative to outcomes that would be attained if everyone followed individualistic best-response reasoning. This advancement of personal interests is made possible—as we will attempt to show in the next section—by a shift in the decision-makers' mode of reasoning from individualistic to reasoning as members of a team.

The second point of our criticism is based on the fact that the use of the average function relies on making interpersonal comparisons of the interacting individuals' payoffs. However, the standard assumptions of the expected utility theory that make numerical representations of individuals' motivations possible by themselves do not allow such comparisons to be made.<sup>10</sup> To illustrate what this means, consider again the Prisoner's Dilemma game presented in

---

<sup>9</sup>See Rawls (1971) and Nozick (1974) for criticism of utilitarianism. Although we use an analogy to their notion of the separateness of persons, our claim here is not based on grounds of moral normativity.

<sup>10</sup>In the literature on the theory of team reasoning this has also been pointed out by Sugden (2000).

Figures 2(A) and 3(A). The numerical representation of the row player’s preferences allows us to say that he or she prefers the outcome  $(C, C)$  to the outcome  $(D, D)$ . It also allows us to say that he or she prefers the outcome  $(C, C)$  to the outcome  $(C, D)$  by a greater extent than he or she prefers the outcome  $(D, C)$  to the outcome  $(C, C)$ . However, it does not allow us to claim that the row player “enjoys” the benefits of the outcome  $(C, C)$  by as much as does the column player.

Although in this paper we do not discuss the technical arguments of why it is so, any payoff function that numerically represents a decision-maker’s preferences or motivations is unique only up to positive affine transformations.<sup>11</sup> What follows from this is that the payoff structure of the game illustrated in Figure 5(B) represents exactly the same motivations of the interacting individuals as the one in Figure 5(A).

	<i>C</i>	<i>D</i>		<i>C</i>	<i>D</i>
<i>C</i>	2, 2	0, 3	<i>C</i>	6, 3	0, 4
<i>D</i>	3, 0	1, 1	<i>D</i>	9, 1	3, 2
	<i>A</i>			<i>B</i>	

Figure 5: The original Prisoner’s Dilemma game (A) and its representation using positive affine transformations<sup>12</sup> of the row and the column players’ payoffs (B)

If the theory of team reasoning makes use of an aggregative payoff function to represent the interests of a team—such as the maximization of the average or the sum of individuals’ personal payoffs—it ceases to be a mere extension of the framework used by the standard rational choice theory and needs to suggest how the required interpersonal comparisons of payoffs are possible. An alternative is to drop the idea that team’s goals are best represented by such aggregative functions and in the next section we propose a representation of team’s interests based on the latter approach. To put this point differently, we hope to show that team play is possible in certain types of interpersonal interactions even without invoking interpersonal comparisons of the interacting players’ payoffs.

<sup>11</sup>This means that if  $u$  is a payoff function representing an individual’s personal motivations, then so is function  $u' = au + c$  where  $a > 0$  and  $c$  are constants. For a detailed discussion of why this is so see, for example, Luce and Raiffa (1957, ch. 2).

<sup>12</sup>The column player’s payoffs in Figure 5(B) were transformed by adding 1 to every number representing his or her payoffs in the original representation of the Prisoner’s Dilemma game in Figure 5(A) while the row player’s payoffs were transformed by multiplying each number by a factor of 3.

## 4 Team’s Interests: A Rank-Based Function

In this section we propose a representation of team’s interests that fits with our intuition about a set of conditions that have to be satisfied in order for individual decision-makers to take part in team play. We start with the proposition of two properties for a candidate function of team’s interests to have:

1. A team’s objective has to be the attainment of an outcome that is beneficial to every member of the team.
2. Decision-makers’ payoffs are not interpersonally comparable, but each individual’s complete preferential ranking of the possible outcomes of a game is commonly known by all players. The derivation of a team’s objective should be done without invoking interpersonal comparisons of players’ payoffs.

The first property captures the idea that there has to be something “in it” for a potential member of a team from team play in order for that individual to participate in the attainment of the team’s goal. As mentioned earlier, we share the intuition behind this property with Sugden (2011, 2015) who suggests that team play has to be recognized as being mutually advantageous by all those partaking in it. Of course, we will need to define what it means for something to be regarded as mutually advantageous as well as how and relative to what that advantage is measured. The second property is a corollary of the standard assumption of non-cooperative game in games of complete information: the common knowledge about the interacting players’ payoffs. It is also in line with the standard axioms of rationality that make numerical representations of decision-makers’ preferences possible but do not entail their interpersonal comparability.

With these properties in mind, we suggest a team’s objective to be the maximal advancement of mutual benefit to the members of the team. As we will show shortly, our proposed function of team’s interests can be operationalized using two types of unit by which the extent of individual and mutual advantage is measured. The choice of unit depends on the type of information conveyed by the payoff structures of games about the interacting decision-makers’ preferences. We start with the case where players’ personal payoff numbers represent merely ordinal preferential rankings of the available outcomes in games. Following this we discuss a scenario in which relative payoff intervals associated with different pairs of outcomes convey meaningful information about players’ relative preferential intensities. In decision-theoretic language, the former refers to ordinal and the latter to cardinal representations of players’ preferences.

There are two ways to motivate our proposed function of team’s interests. One is axiomatic that presents a set of more fine-grained properties for a function of team’s interests to satisfy. The other describes a plausible reasoning process that rational decision-makers may engage in when facing particular types of games. We start with the axiomatic approach first by proposing four axioms to characterize the interests of a team. We will describe a reasoning process that

is in line with these axioms when discussing a number of examples later in this section.

**Axiom 1 (Weak Pareto optimality):** *If an outcome  $y$  of a game is strictly preferred to some other outcome  $x$  by all players, then the outcome  $x$  is not chosen by the team.*

This axiom ensures Pareto efficiency of the selected outcome in a weak sense<sup>13</sup>. We believe this to be an essential feature of any candidate function for representing team’s interests. Since team-reasoning individuals evaluate all outcomes of a game and identify a subset of these as best for a team, it would be odd if they picked an outcome that, from every player’s personal point of view, was worse than some other available alternative.<sup>14</sup>

Before presenting the second axiom we need to introduce a few additional terms. Let the smallest personal payoff that a player can attain from choosing a particular strategy in a game be called that player’s personal security payoff associated with the strategy in question. Recall the Stag Hunt game illustrated in Figure 2(B). The personal security payoff associated with hunting stag for either player is 0, since it is the *minimal* personal payoff that can be attained by playing  $S$ . Similarly, the personal security payoff associated with hunting hare (playing  $H$ ) for either player is 1. Given this and any game, a player can always choose a strategy associated with the highest personal security payoff. This guarantees the player the attainment of a payoff that is at least as high as the security payoff in question, irrespective of what the other players are going to do. In the case of the Stag Hunt game, the *maximal* payoff that any player can guarantee him or her self in this way is 1, which is the personal security payoff associated with hunting hare. This is usually referred to as the player’s personal *maximin* payoff level in the game and the corresponding strategy—the *maximin* strategy.

**Axiom 2 (Preservation of personal security):** *Team play cannot leave any player worse-off than his or her personal maximin payoff level in a game.*

This axiom limits a team’s objective to the attainment of only those outcomes that result in players’ personal payoffs being at least as high as the payoffs that they could secure themselves by playing their *maximin* strategies individually. As such, it defines the lower threshold points in games, below which potential members of a team would, so to speak, not “agree to go” in team play, since they can guarantee themselves a better personal payoff individually. Sugden (2015) uses the same lower threshold point in defining mutually advantageous team

---

<sup>13</sup>An outcome of a game is Pareto efficient in a weak sense if there exists no other outcome in which every player is better off in terms of their personal payoffs. The set of Pareto efficient outcomes in a weak sense is a subset of all Pareto efficient outcomes, since the latter requires there to be no other outcome in which somebody’s payoff could be increased without making anyone else worse off.

<sup>14</sup>Bardsley et al. (2010) consider possible cases where the Pareto criterion may be abandoned in team play. We will briefly return to this in more detail in section 5.

play. (Our approach differs from Sugden's through the imposition of the weak Pareto criterion discussed above. We also extend the notion of mutual benefit by proposing a measure of the extent of individual and mutual advantage presented by different outcomes to the interacting decision-makers, discussed next.)

We now turn to defining individual and mutual advantage and the way these are measured. To do this we introduce the following method for assigning preferential rank values to the available outcomes in games based on each player's personal preferential ordering of those outcomes. For a particular player, the least preferred outcome in a game is assigned the preferential rank value 0. The second least preferred outcome is assigned the preferential rank value 1 and so on. A shift from one outcome to another that results in an increase of the assigned preferential rank value by 1 for some player is said to advance that player's personal interests by 1 unit and, more generally, an increase of the assigned preferential rank value by  $k$  is said to advance that player's personal interests by  $k$  units.

**Individual advantage:** An outcome of a game is individually advantageous to a particular player if that player preferentially ranks this outcome above the outcome(s) associated with his or her personal *maximin* payoff level in the game. The extent of individual advantage provided by an outcome to a particular player is given by the number of units this outcome advances that player's personal interests relative to the outcome(s) associated with his or her personal *maximin* payoff level in the game.

**Mutual advantage:** An outcome of a game is mutually advantageous to the interacting players if each player preferentially ranks this outcome above the outcome(s) associated with his or her personal *maximin* payoff level in the game. The extent of mutual advantage provided by an outcome to the interacting decision-makers is given by the number of units this outcome advances all players' personal interests in parallel relative to the outcome(s) associated with each player's personal *maximin* payoff level in the game.

For an example imagine a two-player game where some outcome  $x$  advances the first and the second player's personal interests relative to the outcomes associated with their personal *maximin* payoff levels in the game by 1 and 2 units respectively. Since this outcome advances the two players' personal interests in parallel by 1 unit, it is said to provide 1 unit of mutual advantage to the interacting decision-makers. The additional unit of advancement of the second player's personal interests provided by this outcome represents individual advantage to the second player over and above the 1 unit of mutual advantage. We will illustrate this with a number of more concrete examples later in this section.

**Axiom 3 (Maximal mutual advantage):** *An outcome selected by a team has to be maximally mutually advantageous.*

This axiom has two implications. First, if a game contains a mutually advantageous outcome, then any outcome that is selected by a team must also be mutually advantageous. Second, if a game contains multiple mutually advantageous outcomes, then the outcome(s) selected by the team must provide *maximal* mutual advantage. In other words, if there is an outcome that provides more mutual benefit to the interacting decision-makers than some outcome  $x$ , then  $x$  is not chosen by the team.

Given the three axioms above we propose the following function to represent the interests of a team.

**Rank-based function of team’s interests:** Maximize the minimum number of units by which individuals’ personal interests are advanced among the interacting players relative to each player’s threshold point—the outcome(s) associated with his or her *maximin* payoff level in the game.

It can be shown that this function satisfies the three axioms introduced earlier (for proofs see Appendix A). It is not, however, the only function that does so. Consider a two-player game where a particular outcome  $x$  advances each of the two players’ personal interests by 2 units while some other outcome  $y$  advances the first and the second player’s personal interests by 2 and 3 units respectively. The extent of mutual advantage provided by either of the two outcomes is the same—2 units. Assuming there to be no outcomes providing a higher extent of mutual advantage in this game, the rank-based function of team’s interests selects both  $x$  and  $y$ . As such, any function that would further discriminate between these two outcomes would also satisfy the three axioms introduced above. However, any such additional discrimination would go beyond the notion of mere maximization of the extent of mutual advantage attained by the interacting players. We therefore believe that the proposed function best captures the interacting individuals’ primary motivation to satisfy their personal preferences in a mutually beneficial way without invoking any additional motivational attitudes, such as benevolence towards others in team play. We formalize the latter idea using Axiom 4 below.

**Axiom 4 (Indifference between equal extent of mutual advantage):** *If some outcome  $x$  is in a team’s choice set<sup>15</sup> and some other outcome  $y$  provides the same extent of mutual advantage as  $x$ , then  $y$  is in the team’s choice set as well.*

This says that all that matters in determining the optimal outcomes for a team is the extent of mutual advantage that those outcomes provide to the interacting individuals. This axiom ensures uniqueness of the proposed function of team’s interests (for proof see Appendix B).

There is a connection between the above function of team’s interests and the idea of rational cooperation discussed by Gauthier (2013). For Gauthier, rational cooperation is the attainment of Pareto efficiency through the maximization

---

<sup>15</sup>To say that an outcome is in a team’s choice set means the same as to say that it is selected as one of the optimal outcomes for the team.

of the minimum level of individual benefit across individuals similar to the *maximin* principle in mutually advantageous team play described above. Although Gauthier mentions relative threshold points, below which decision-makers would refuse to cooperate, he does not provide a clear characterization of what those threshold points are and how they are derived. Also, Gauthier's justification of the proposed function of rational cooperation (the idea of *maximin* proportionate gain in his text) is not entirely clear and the axiomatic approach presented here is a possible way for filling these gaps.

#### 4.1 Formalization

Let  $\Gamma$  be a normal form game defined as a triple  $(I, S_i, u_i)$  where  $I = \{1, 2, \dots, m\}$  is a finite set of  $m$  players,  $S_i$  is a set of pure strategies available to player  $i \in I$ , and  $u_i(\mathbf{s})$  is a payoff function that assigns to every player  $i \in I$  a personal payoff for each outcome in the game. An outcome of a game is defined as a strategy profile  $\mathbf{s} = (s_1, \dots, s_m)$  where  $s_i \in S_i$  is a particular pure strategy chosen by player  $i \in I$ . Let  $\Sigma$  be the set of all possible strategy profiles (i.e. outcomes) in the game  $\Gamma$  and  $\mathbf{s}^{\vee i} \in \Sigma$  be a strategy profile that is associated with player  $i$ 's *maximin* payoff level in the game.

Each player  $i \in I$  has a personal preferential ranking of all the strategy profiles in the set  $\Sigma$ . Let  $\{x \in \mathbb{Z}^* \mid x \leq k\}$  be a set of non-negative integers from 0 to  $k$  where  $k \leq n - 1$  and  $n$  is the total number of strategy profiles available in the game  $\Gamma$ . Let  $\wp_i : \Sigma \rightarrow \{x \in \mathbb{Z}^* \mid x \leq k\}$  be a ranking function that, for each player  $i \in I$ , assigns a personal preferential rank value to each strategy profile in the set  $\Sigma$  as follows: the strategy profile that player  $i \in I$  prefers the least in the game  $\Gamma$  is assigned the preferential rank value 0, the second least preferred strategy profile is assigned the preferential rank value 1 and so on.

The rank-based function of team's interests  $\mathcal{F}^\tau : \mathcal{P}(\Sigma) \rightarrow \mathcal{P}(\Sigma)$ , where  $\mathcal{F}^\tau(\Sigma) = \Sigma^\tau$  and  $\Sigma^\tau \subseteq \Sigma$ , is a choice function that selects a subset from the set of all possible strategy profiles of the game  $\Gamma$  such that each selected strategy profile maximizes the minimum difference, across all players, between the preferential rank value of the selected profile and the preferential rank value of the strategy profile associated with a particular player's personal *maximin* payoff level in the game. In other words, each element  $\mathbf{s}^\tau \in \Sigma^\tau$  is such that

$$\mathbf{s}^\tau \in \arg \max_{\mathbf{s} \in \Sigma} \left\{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \right\} :=$$

$$\left\{ \mathbf{s}^\tau \mid \forall \mathbf{s} \in \Sigma : \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \leq \min_{i \in I} [\wp_i(\mathbf{s}^\tau) - \wp_i(\mathbf{s}^{\vee i})] \right\}$$

#### 4.2 Payoff Intervals in Individual and Mutual Advantage

The operationalization of the rank-based function of team's interests presented thus far is based on the interacting players' ordinal preferential rankings of outcomes in games and it ignores the relative preferential intensities in the players' pairwise comparisons of those outcomes. Suppose a player prefers some outcome

$x$  to the outcome  $y$  and the outcome  $y$  to the outcome  $z$ . The operationalization of the function presented above considers the player's preferential ranking of the three outcomes, but ignores the information on whether the player prefers  $x$  to  $y$  by a greater, lesser or the same extent as he or she prefers  $y$  to  $z$ . The standard axioms of rationality, however, make the ratios of payoff intervals between different pairs of outcomes to be meaningful representations of decision-makers' relative preferential intensities. We now present a slight modification of the above approach that takes the information about such preferential intensities into account.

We introduce the following modification for determining the number of units by which a shift from one outcome to another is said to advance a particular player's personal interests. First, each player's personal payoffs associated with the available outcomes in games are normalized to assign the least preferred outcome the payoff value of 0 and the most preferred outcome the payoff value of 100. This is done by applying an appropriate positive affine transformation of personal payoffs. For example, if a particular player's preferences over four outcomes in a game are represented by payoff values 0, 1, 2 and 3, the normalization transforms these into 0,  $33\frac{1}{3}$ ,  $66\frac{2}{3}$  and 100 respectively. After this normalization, a shift from one outcome to another that results in an increase in the assigned normalized personal payoff value of 1 for some player is said to advance that player's personal interests by 1 unit. Apart from this modification, the three axioms introduced earlier, the definitions of individual and mutual advantage and the description of the rank-based function of team's interests all remain unchanged.

In the formal representation of the modified approach  $\Gamma$  is a normal form game defined as a triple  $(I, S_i, u_i^*)$  where  $u_i^*(\mathbf{s})$  is a normalized payoff function that assigns to every player  $i \in I$  a personal payoff for each outcome in the game in a way that the least and the most preferred outcomes are given the values 0 and 100 respectively. Each element  $\mathbf{s}^\tau \in \Sigma^\tau$  is now such that

$$\mathbf{s}^\tau \in \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [u_i^*(\mathbf{s}) - u_i^*(\mathbf{s}^{\vee i})] \} := \\ \{ \mathbf{s}^\tau \mid \forall \mathbf{s} \in \Sigma : \min_{i \in I} [u_i^*(\mathbf{s}) - u_i^*(\mathbf{s}^{\vee i})] \leq \min_{i \in I} [u_i^*(\mathbf{s}^\tau) - u_i^*(\mathbf{s}^{\vee i})] \}$$

### 4.3 Examples

We now turn to discussing a number of examples. First, we apply the function of team's interests to four well known two-player games with multiple Nash equilibria. We then revisit the Prisoner's Dilemma—a case involving a unique Nash equilibrium. Since the prescriptions of the rank-based function are the same for both types of unit used to measure the extent of individual and mutual advantage in all these cases, we will focus solely on the ordinal representations of players' preferences where the unit in question is the preferential rank value described on page 13. We will end this section by presenting an example where prescriptions of the rank-based function diverge for the two types of unit—the preferential rank value and the normalized payoff.



### 4.3.1 The Hi-Lo

Recall the Hi-Lo game illustrated in Figure 1. The row and the column players' personal preferential rankings of the four outcomes are shown below with numbers representing the corresponding preferential rank values and arrows indicating outcomes associated with each player's *maximin* payoff level in the game.

$$\text{Row player: } \left[ \begin{array}{cc} (Hi, Hi) & 2 \\ (Lo, Lo) & 1 \\ \Rightarrow (Hi, Lo)(Lo, Hi) & 0 \end{array} \right] \text{ Column player: } \left[ \begin{array}{cc} (Hi, Hi) & 2 \\ (Lo, Lo) & 1 \\ \Rightarrow (Hi, Lo)(Lo, Hi) & 0 \end{array} \right]$$

Relative to the *maximin* payoff level, the outcome  $(Hi, Hi)$  advances each player's personal interests by 2 units. The outcome  $(Lo, Lo)$  does so by 1 unit. Hence, the rank-based function of team's interests identifies  $(Hi, Hi)$  as the best outcome for a team and its attainment as the team's objective. As a result, for any individual who reasons as a member of a team,  $(Hi, Hi)$  is the only rational outcome in this game.

A plausible reasoning process by which individual players may arrive at this conclusion can be described as follows. Having worked out all the best-response strategies in this game, individualistically reasoning decision-makers face a Nash equilibrium selection problem. This leaves them stuck with no further indication of what actions they ought to perform independently from each other in order to best promote their personal interests. The next question they ask themselves is "what would be best for both of us in this situation?". As soon as they do so, they start reasoning as members of a team and identify the uniquely rational outcome from the perspective of team's interests.

According to the deliberative process just described, the adoption of the team mode of reasoning comes about as a result of extensive thinking about the game by otherwise individualistically rational agents who are able to work out each other's best-response strategies and realize that these leave them with a Nash equilibrium selection problem. As such, it uses as its starting point the standard assumption of orthodox game theory that the basic mode of reasoning used by the interacting decision-makers is individualistic best-responding.<sup>16</sup>

### 4.3.2 The Stag Hunt

Recall the Stag Hunt game illustrated in Figure 2(B). The row and the column players' preferential rankings of outcomes are shown below.

$$\text{Row player: } \left[ \begin{array}{cc} (S, S) & 2 \\ \Rightarrow (H, S)(H, H) & 1 \\ (S, H) & 0 \end{array} \right] \text{ Column player: } \left[ \begin{array}{cc} (S, S) & 2 \\ \Rightarrow (S, H)(H, H) & 1 \\ (H, S) & 0 \end{array} \right]$$

Since either player can obtain a payoff of 1 by hunting hare (playing  $H$ ) irrespective of what the other player is going to do, the personal *maximin* payoff

<sup>16</sup>For other possible suggestions of what may trigger shifts in individuals' mode of reasoning, some of which abandon this assumption, see, for example, Bacharach (2006), Sugden (2003) or Gold and Sugden (2007a,b).

level is associated with the pair of outcomes  $(H, S)$  and  $(H, H)$  for the row, and  $(S, H)$  and  $(H, H)$  for the column players. Since  $(S, S)$  is the only mutually advantageous outcome relative to these threshold points, the rank-based function of team's interests identifies it as the uniquely rational solution of the game for players who reason as members of a team.

With regards to this game in particular, it is important to mention other solution concepts to the Nash equilibrium selection problem that are based on risk or uncertainty aversion with respect to choosing a particular action. Notice that in the Stag Hunt game it can be argued that the Nash equilibrium  $(H, H)$  is in some sense safer than the equilibrium  $(S, S)$ . This is because either player is guaranteed a certain payoff from hunting hare (playing  $H$ ) whereas the high payoff from hunting stag (playing  $S$ ) crucially depends on the other player's choice. While we don't review the various models of risk and uncertainty aversion in this paper (some of which are based on probabilistic assessments of players' actions while others on the sizes of foregone payoffs in cases of deviations from equilibrium play) we do believe that these considerations may play a role in determining the likelihood of individuals' adoption of the team mode of reasoning. However, since existing experimental data suggests that in one-shot interactions people tend to play  $S$  slightly more than 60% of the time, we do suggest that a shift in individuals' mode of reasoning to that as members of a team in one-shot versions of this game is likely.<sup>17</sup>

### 4.3.3 The Chicken

Recall the Chicken game illustrated in Figure 4(B). The row and the column players' preferential rankings of outcomes are shown below.

$$\text{Row player: } \left[ \begin{array}{l} (U, L) \quad 3 \\ (D, L) \quad 2 \\ (D, R) \quad 1 \\ (U, R) \quad 0 \end{array} \right] \Rightarrow \text{Column player: } \left[ \begin{array}{l} (D, R) \quad 3 \\ (D, L) \quad 2 \\ (U, L) \quad 1 \\ (U, R) \quad 0 \end{array} \right]$$

Here, similarly as in the Stag Hunt game above,  $(D, L)$  is the only mutually advantageous outcome relative to the two players' threshold points — their personal *maximin* payoff levels in the game. As such, the rank-based function of team's interests identifies its attainment as the objective for the team.

The deliberative process by which individual decision-makers may arrive at this conclusion is described as follows. Thinking individualistically, both players identify the two Nash equilibria in this game:  $(U, L)$  and  $(D, R)$ . From a personal point of view, the row player prefers the attainment of the outcome  $(U, L)$ . At the same time he or she recognizes the column player's preference for the attainment of the outcome  $(D, R)$ . If both players were to pursue their

<sup>17</sup>Rankin et al. (2000) report an experiment in which different versions of the Stag Hunt game are played repeatedly, but in a way that aims to induce one-shot reasoning of participating individuals in each round (by changing labels associated with the available actions and varying sizes of payoffs resulting from different outcomes each time the game is played). In this setting virtually everyone switches to playing  $S$  over time, which seems to support our suggestion.

preferred options, they would end up with the outcome  $(U, R)$ , which is the worst of all possible outcomes for both. The Nash equilibrium selection problem, as previously, leaves them stuck with no further indication of what actions they ought to perform independently from each other in order to best promote their personal interests. At this stage they ask themselves “what would be best for both of us in this situation?”. As soon as they do so, they identify the outcome  $(D, L)$  as the uniquely rational solution of this game from the perspective of team’s interests. Now, since the outcome  $(D, L)$  is not a Nash equilibrium itself, they may each consider unilateral deviation from team play. However, the previously recognized conflict of their personal interests associated with the two Nash equilibria prevents them from doing so, emphasizing the outcome  $(D, L)$  as the only mutually advantageous solution of the game.

#### 4.3.4 The Divide-the-Cake

Another interesting case involving multiple Nash equilibria is the Divide-the-Cake game illustrated in Figure 6, which is a particularly simple version of the well known Nash Bargaining game. In this game two players are presented

	$0$	$1$	$2$	$3$	$4$
$0$	0, 0	0, 1	0, 2	0, 3	0, 4
$1$	1, 0	1, 1	1, 2	1, 3	0, 0
$2$	2, 0	2, 1	2, 2	0, 0	0, 0
$3$	3, 0	3, 1	0, 0	0, 0	0, 0
$4$	4, 0	0, 0	0, 0	0, 0	0, 0

Figure 6: The Dividing of a Cake game

with a cake that is cut into four equal-sized pieces and simultaneously place a demand for the number of pieces for themselves (from  $0$  to  $4$ ). If the sum of their demanded pieces does not exceed  $4$ , they both get what they asked for. If, on the other hand, the sum exceeds  $4$ , they both get nothing. The game has six Nash equilibria:  $(4, 0)$ ,  $(3, 1)$ ,  $(2, 2)$ ,  $(1, 3)$ ,  $(0, 4)$  and an inefficient  $(4, 4)$ . The row and the column players’ personal rankings of outcomes are shown below.<sup>18</sup>

<sup>18</sup>There are 11 outcomes in total that yield a payoff of 0 to both players. For the sake of brevity they are all represented by the outcome  $(0, 0)$  in the two rankings.

$$\begin{array}{l}
\text{Row player:} \\
\Rightarrow
\end{array}
\left[ \begin{array}{cc}
(4, 0) & 4 \\
(3, 0)(3, 1) & 3 \\
(2, 0)(2, 1)(2, 2) & 2 \\
(1, 0)(1, 1)(1, 2)(1, 3) & 1 \\
(0, 0)(0, 1)(0, 2)(0, 3)(0, 4) & 0
\end{array} \right]$$

$$\begin{array}{l}
\text{Column player:} \\
\Rightarrow
\end{array}
\left[ \begin{array}{cc}
(0, 4) & 4 \\
(0, 3)(1, 3) & 3 \\
(0, 2)(1, 2)(2, 2) & 2 \\
(0, 1)(1, 1)(2, 1)(3, 1) & 1 \\
(0, 0)(1, 0)(2, 0)(3, 0)(4, 0) & 0
\end{array} \right]$$

Here, relative to the two individuals' *maximin* payoff levels, the outcome  $(2, 2)$  advances each player's personal interests by 2 units. Since every other outcome either advances one of the players' personal interests by only 1 unit or is not mutually advantageous (i.e. does not advance one of the players' personal interests at all) the rank-based function of team's interests identifies the team's objective with the attainment of the outcome  $(2, 2)$ . This usually appeals to most decision-makers and is supported by experimental results.<sup>19</sup>

The above result is in line with Nash's bargaining solution of this game.<sup>20</sup> Misyak and Chater (2014) propose a theory of virtual bargaining as an alternative mode of reasoning that individuals may adopt when choosing among the available outcomes and courses of actions in various types of games. According to the proposed model the interacting decision-makers are said to be undergoing implicit mental bargaining processes in an attempt to work out possible agreeable outcomes in a somewhat similar fashion as it is done with the rank-based function of team's interests presented here. There are, however, some differences between the two approaches.

First, the existing theories of bargaining generally rely on the existence of a unique reference outcome that obtains when individuals fail to reach an agreement following a bargaining process. In the Divide-the-Cake game this is assumed to be the outcome in which both players gain nothing. In the Prisoner's Dilemma game this is usually the Nash equilibrium  $(D, D)$ —the outcome that both players can fall back to in case of a failure of reaching an agreement with regards to anything else. In most other cases that we address in this paper, however, there is no such unique reference point. This does not pose a problem for the proposed rank-based function of team's interests, since its outputs are determined not in relation to a unique reference outcome but by considering multiple threshold points—one for each player in a game.

<sup>19</sup>See Nydegger and Owen (1974) for an experiment in which two players are asked to divide \$1 among themselves and virtually everybody agrees on a 50%-50% split. Note, however, that for an odd number of slices the prescription of the rank-based function will not yield a unique outcome. For example, in case of five slices the function selects three outcomes as optimal for the team:  $(2, 2)$ ,  $(3, 2)$  and  $(2, 3)$ . We will discuss the indeterminacy of the rank-based function in more detail later in the text.

<sup>20</sup>The fact that the predicted outcomes are the same in this particular example is a coincidence and the proposed rank-based function of team's interests is more in line with the bargaining solution presented by Kalai and Smorodinsky (1975). Although the differences and implications of the two bargaining models are interesting, we do not discuss them in more detail here. For a discussion of Nash's bargaining model see, for example, Luce and Raiffa (1957, ch. 6).

The second difference lies with the fact that Nash’s bargaining model (which is presented as a possible starting point in the development of the theory of virtual bargaining) does not entirely avoid the need to make interpersonal comparisons of individuals’ payoffs. This is because Nash’s solution requires the interacting decision-makers’ payoff functions to be unique only up to positive linear transformations and not positive affine transformations as we assume here. This difference, however, can be overcome by adopting a different bargaining model, such as the one presented by Kalai and Smorodinsky (1975).<sup>21</sup>

### 4.3.5 The Prisoner’s Dilemma

Recall the original and the amended Prisoner’s Dilemma games illustrated in Figures 2(A) and 4(A) respectively. The row and the column players’ preferential rankings of the four outcomes and their personal *maximin* threshold points (which are the same in both games) are shown below.

$$\text{Row player: } \left[ \begin{array}{l} (D, C) \ 3 \\ (C, C) \ 2 \\ (D, D) \ 1 \\ (C, D) \ 0 \end{array} \right] \Rightarrow \quad \text{Column player: } \left[ \begin{array}{l} (C, D) \ 3 \\ (C, C) \ 2 \\ (D, D) \ 1 \\ (D, C) \ 0 \end{array} \right]$$

The outcomes  $(C, D)$  and  $(D, C)$  lie below one or the other player’s threshold point. Hence their attainment is not viable in team play. Since  $(C, C)$  is the only outcome that advances both players’ personal interests relative to their personal *maximin* payoff levels, it is identified as the unique solution of the game from the perspective of team’s interests.

Even though this game has a unique Nash equilibrium, the motivation to opt for a mutually beneficial outcome may prevail. This is because players are aware of the fact that, if each of them pursues his or her personal goals individually, they will end up at a suboptimal, Pareto inefficient outcome (in the weak sense of Pareto efficiency). According to orthodox game theorists, games such as the Prisoner’s Dilemma are trivial, since they contain a unique rational outcome in terms of individualistic reasoning. Many people outside this field, however, feel that cooperation in the Prisoner’s Dilemma game is not unreasonable. We believe that the major source of controversy and disagreement about this game may lie with the fact that some people tend to think about it from the perspective of team’s interests.

It is possible to provide two interpretations of what happens when decision-makers’ mode of reasoning undergoes a shift from individualistic to reasoning as members of a team in this case. According to one interpretation, individuals who reason as members of a team identify the uniquely rational outcome from the perspective of team’s interests with other options no longer appearing to

---

<sup>21</sup>The underlying mechanism of the rank-based function of team’s interests is somewhat similar to that behind the egalitarian solution to bargaining problems presented by Kalai (1977) and Myerson (1977). The egalitarian solution is in line with the maximization of the minimum advancement of players’ personal payoffs relative to a given reference point. However, since it is concerned with equal advancement of players’ payoffs, the model is also reliant on the interpersonal comparability of those payoffs.

them as rational solutions of this game. According to another interpretation and one that is in line with the deliberative processes we discussed in the preceding examples, individuals who reason as members of a team recognize the existence of two rational solutions—one in terms of individualistic reasoning and the other from the perspective of team’s interests. This leaves their chosen courses of actions undetermined and rationalizable in two ways. It is only the latter interpretation, however, that turns these games into genuine dilemmas for the interacting decision-makers, since only those individuals who will identify in such cases two distinct, differently rationalizable solutions will be puzzled about how to proceed.

#### 4.3.6 Preferential Rank Value vs. Normalized Payoff

In all examples discussed thus far the prescriptions of the rank-based function are same for the two types of unit that can be used to measure the extent of individual and mutual advantage provided by the different outcomes in games. For an example of where these prescriptions differ consider the game illustrated in Figure 7(A). This game has a unique Nash equilibrium:  $(D, R)$ . The row

	<i>L</i>	<i>R</i>		<i>L</i>	<i>R</i>
<i>U</i>	9, 3	0, 2	<i>U</i>	90, 6	0, 4
<i>M</i>	8, 50	0, 2	<i>M</i>	80, 100	0, 4
<i>D</i>	10, 0	1, 1	<i>D</i>	100, 0	10, 2
		<i>A</i>			<i>B</i>

Figure 7: An example with different prescriptions for the two types of unit used to measure the extent of individual and mutual advantage

and the column players’ preferential rankings of outcomes are shown below with each player’s normalized payoffs given in parentheses. The same game using the two players’ normalized payoffs is illustrated in Figure 7(B).

$$\begin{array}{l}
 \text{Row player:} \\
 \Rightarrow \\
 \text{Column player:} \\
 \Rightarrow
 \end{array}
 \left[ \begin{array}{l}
 (D, L) \quad 4 \quad (100) \\
 (U, L) \quad 3 \quad (90) \\
 (M, L) \quad 2 \quad (80) \\
 (D, R) \quad 1 \quad (10) \\
 (U, R)(M, R) \quad 0 \quad (0) \\
 (M, L) \quad 4 \quad (100) \\
 (U, L) \quad 3 \quad (6) \\
 (U, R)(M, R) \quad 2 \quad (4) \\
 (D, R) \quad 1 \quad (2) \\
 (D, L) \quad 0 \quad (0)
 \end{array} \right]$$

The two mutually advantageous outcomes are  $(U, L)$  and  $(M, L)$ . Using the preferential rank value numbers 0 to 4, the outcome  $(U, L)$  advances both players' personal interests by 2 units. The outcome  $(M, L)$  advances the row and the column players' personal interests by 1 and 3 units respectively. The extent of mutual advantage provided by the outcome  $(U, L)$  is 2 and that provided by the outcome  $(M, L)$  is 1. As such, the approach that considers preferential rank value numbers to measure the extent of individual and mutual advantage selects the outcome  $(U, L)$  as optimal for the team.

It may seem, however, that the outcome  $(M, L)$  is in some sense better for the team than the outcome  $(U, L)$ . This is because the former results in the row player attaining a payoff that is almost as good as the best possible payoff to him or her in this game and the column player attaining a payoff that is far better than any of the other five possibilities. This reasoning, however, is based on the idea that payoff intervals between different pairs of outcomes convey meaningful information about the players' preferential intensities and can be captured using normalized payoffs to measure the extent of mutual advantage provided by the two outcomes. Using the normalized payoff values 0 to 100, the outcome  $(U, L)$  advances the row and the column players' personal interests by 80 and 4 units respectively. The outcome  $(M, L)$  does so by 70 and 98 units respectively. The *minimal* advancement of personal interests across players as well as the extent of mutual advantage provided by the outcome  $(U, L)$  is now 4 and that provided by the outcome  $(M, L)$ —70. As a result, the rank-based function of team's interests operationalized using normalized payoffs favours the outcome  $(M, L)$ .

#### 4.4 Indeterminacy of the Rank-Based Function

It is not always the case that the proposed rank-based function yields a unique solution to a group of individuals who are reasoning as members of a team. For an example consider a version of the Chicken game illustrated in Figure 8. As in the previous case, this game has two Nash equilibria:  $(U, L)$  and  $(D,$

	$L$	$R$
$U$	3, 2	0, 0
$D$	1, 1	2, 3

Figure 8: The Chicken game (version 2)

$R)$ . The row and the column players' preferential rankings of outcomes and their personal *maximin* threshold points are shown below with each player's normalized payoffs given in parentheses.

$$\text{Row player: } \left[ \begin{array}{l} (U, L) \quad 3 \quad (100) \\ (D, R) \quad 2 \quad (66\frac{2}{3}) \\ \Rightarrow \quad (D, L) \quad 1 \quad (33\frac{1}{3}) \\ (U, R) \quad 0 \quad (0) \end{array} \right] \quad \text{Column player: } \left[ \begin{array}{l} (D, R) \quad 3 \quad (100) \\ (U, L) \quad 2 \quad (66\frac{2}{3}) \\ \Rightarrow \quad (D, L) \quad 1 \quad (33\frac{1}{3}) \\ (U, R) \quad 0 \quad (0) \end{array} \right]$$

Using the preferential rank value numbers 0 to 3, the outcome  $(U, L)$  advances the row and the column players' personal interests by 2 and 1 units respectively. The outcome  $(D, R)$  does so by 1 and 2 units respectively. The *minimal* advancement of personal interests across the two players as well as the extent of mutual advantage that is provided by both outcomes is 1. Hence the attainment of either one of these two outcomes remains a viable goal for a team. (The same result obtains using the normalized payoff values 0 to 100.) As a result, this game poses a coordination problem for individualistically reasoning decision-makers as well as those who reason as members of a team and further methods for resolving this game may be sought. We will briefly return to this issue in the next section.

#### 4.5 Normative vs. Descriptive Status

It may be asked whether the theory of team reasoning operationalized using the proposed function of team's interests is a normative or a descriptive theory of choice in interdependent decision situations. The answer, in our view, is somewhat mixed and depends largely on the type of decision problem individuals face. We propose to classify games using the following four categories: [1] games with multiple Nash equilibria, [2] games with a unique inefficient Nash equilibrium (in the weak sense of Pareto efficiency) [3] games with a unique efficient Nash equilibrium (in the weak sense of Pareto efficiency) and [4] games with no Nash equilibria, where in all four categories we refer to Nash equilibria in pure strategies alone.

In the case of category [1] — games with multiple Nash equilibria — we suggest that the theory of team reasoning has a strong normative appeal and that a shift in the interacting players' mode of reasoning from individualistic to that as members of a team is likely from the descriptive point of view as well. Although individualistic best-response reasoning provides a number of rational solutions, it does not resolve these games fully in terms of providing a definitive course of actions for the decision-makers to take. The theory of team reasoning operationalized as above, however, suggests that best-response reasoning is not the endpoint of rational deliberation and provides a rational resolution of these games based on the notion of mutually advantageous advancement of the interacting players' personal interests. As such, if one agrees with the proposed definition of mutual advantage and the suggested axioms to characterize the interests of a team, the rank-based function of team's interests provides a normatively compelling basis for resolving these games. Whether the switch to reasoning as a member of a team fully resolves such games, of course, depends on whether the outcome selected by the function of team's interests is unique. This, as illustrated earlier, may not always be the case and other methods for resolving such scenarios may be required. One possibility is a search for differ-



entiating features of the contemplated outcomes in order to single out one of them for the sake of successful coordination of players' actions. For example, if the rank-based function identified three outcomes yielding payoff pairs  $(2, 2)$ ,  $(2, 3)$  and  $(3, 2)$  as equally good for a team of two players<sup>22</sup>, the outcome  $(2, 2)$  may be singled out due to its uniqueness. In many cases, however, a switch to reasoning as a member of a team does resolve these games definitively and we take this switch to be a natural progression in a player's rational deliberative process in an attempt to resolve such games when best-response reasoning fails to do so.

In the case of category [2] the answer is less clear. These games have a unique Nash equilibrium and, as such, individualistic best-response reasoning resolves them fully. However, having identified the uniquely optimal outcomes based on individualistic best-response reasoning, the interacting players may easily recognize their unappealing inefficiencies. This, in turn, may prompt them to identify mutually advantageous outcomes from the perspective of team's interests. From the descriptive point of view we expect this to happen often, making such cases genuine dilemmas for the interacting decision-makers due to the existence of multiple differently rationalizable resolutions of these games. From the normative point of view, however, multiple courses of actions remain possible: some based on the principles of best-response reasoning and others based on those of mutually advantageous team play. As a result, a decision-maker who has recognized both modes of reasoning in such games will have to decide on which underlying principles to base his or her choice.

In the case of category [3] individualistic best-response reasoning resolves these games fully and efficiently. As such, there is no need for team reasoning to resolve these games efficiently and we do not expect people to reason as members of a team from the descriptive point of view either.

Finally, in the case of category [4], similarly as was the case with category [1], we suggest the theory of team reasoning to have a strong normative appeal, since individualistic best-response reasoning provides no definitive solutions in these games. In particular scenarios—those in which the function of team's interests is able to identify uniquely rational solutions—we also expect a shift in decision-makers' mode of reasoning to provide a descriptively accurate explanation of the chosen courses of action.

## 5 Conclusion

In this paper we argued against the operationalization of the theory of team reasoning with the use of the maximization of the average of individuals' personal payoffs as a representation of team's interests. Our arguments focused on its reliance on making interpersonal comparisons of the interacting players' payoffs and a possible advocacy of a complete sacrifice of some individuals' personal interests for the benefit of others. While only a subset of the existing texts on the theory of team reasoning make a reference to the maximization of average

---

<sup>22</sup>See earlier footnote 19 for an example where this may be the case.

payoffs, most of the remaining literature in this field does not endorse a specific function and instead focuses on specifying a number of general conditions for a candidate function of team’s interests to satisfy. We attempted to fill this gap by proposing a rank-based function of team’s interests as a possible alternative that is in line with the orthodox conception of payoffs and is based on the notion of mutual advantage in team play that is compatible with the idea suggested by Sugden (2011, 2015). We extended the notion of mutual benefit by presenting a measure of the extent of individual and mutual advantage provided by the available outcomes to the interacting players in games and axiomatized the formal representation of team’s interests.

While the proposed rank-based function of team’s interests fits with some experimental findings from games discussed in this paper, further empirical tests will need to be constructed to test this model’s empirical validity. In principle this task is possible, since the rank-based function of team’s interests provides testable predictions in many games.

The second area requiring further research concerns cases where the proposed rank-based function of team’s interests does not yield a unique solution. Bardsley et al. (2010) consider possible scenarios where the Pareto efficiency criterion may be abandoned in such cases. An example is the Enlarged Hi-Lo game illustrated in Figure 9. In this game the rank-based function identifies team’s

	<i>Hi 1</i>	<i>Hi 2</i>	<i>Lo</i>
<i>Hi 1</i>	10, 10	0, 0	0, 0
<i>Hi 2</i>	0, 0	10, 10	0, 0
<i>Lo</i>	0, 0	0, 0	9, 9

Figure 9: The Enlarged Hi-Lo game

interests with the attainment of outcomes (*Hi 1*, *Hi 1*) or (*Hi 2*, *Hi 2*) but does not discriminate further among the two. One way to proceed for a team in this case is to abandon the Pareto efficiency criterion and instead focus on the attainment of the outcome (*Lo*, *Lo*) for the sake of successful coordination of players’ actions. This suggests that the function of team’s interests may need to be developed further to account for such scenarios. Another possibility, however, is to separate the question of which outcomes in a considered game are the best from the point of view of mutually advantageous team play from the question of how to coordinate players’ actions having identified the best outcomes. From this point of view, the Enlarged Hi-Lo game could be seen to have three rational solutions in terms of individualistic reasoning and two rational solutions in terms of mutually advantageous team play. This could prompt the decision-makers to seek further methods for coordinating their actions among the team-optimal

outcomes first and, if no coordination was possible among these, revert back to considering outcomes that were suboptimal from the perspective of team’s interests—such as the outcome  $(Lo, Lo)$  in the above example—to seek possible coordination among those.

Finally, the exact ways by which the interacting individuals arrive at the conclusions predicted by the rank-based function of team’s interests depend on what is responsible for triggering shifts in people’s mode of reasoning from individualistic to reasoning as members of a team and vice versa. While we did not discuss the various accounts addressing this question in the existing literature, we presented a possible suggestion according to which a shift in individuals’ mode of reasoning may come about as a result of extended reasoning in games by otherwise individualistically rational agents. Further research into this possibility would allow to bridge the gap between the two modes of reasoning within the framework of non-cooperative game theory.

## Appendix A

We here show that the rank-based function of team’s interests presented above satisfies the four axioms introduced in section 4. Proofs presented below are tailored to the case where the extent of individual and mutual advantage provided by different outcomes to the interacting players in games is measured using preferential rank values discussed on page 13. For a version of these proofs where the extent of both types of advantage is measured using normalized payoffs discussed on page 16, all instances of  $\wp_i(\cdot)$  need to be replaced with  $u_i^*(\cdot)$ .

**Axiom 1 (Weak Pareto optimality):** *If an outcome  $y$  of a game is strictly preferred to some other outcome  $x$  by all players, then the outcome  $x$  is not chosen by the team.*

Suppose that an outcome selected by the rank-based function of team’s interests does not satisfy this axiom. In other words, suppose that some outcome  $x$  is selected by a team when some other outcome  $y$  is strictly preferred to  $x$  by all players. Following the notation introduced in section 4.1, let the two outcomes  $x$  and  $y$  be defined as strategy profiles  $\mathbf{s}^x, \mathbf{s}^y \in \Sigma$  respectively. Since  $\mathbf{s}^y$  is strictly preferred to  $\mathbf{s}^x$  by every player  $i \in I$ , it is the case that (i) for every player  $i \in I$ ,  $\wp_i(\mathbf{s}^y) > \wp_i(\mathbf{s}^x)$ , where  $\wp_i(\mathbf{s})$  is player  $i$ ’s preferential rank value associated with the strategy profile  $\mathbf{s} \in \Sigma$ .

By Axiom 2, any outcome that is selected by a team is, for every player, at least as good as the outcome associated with his or her personal *maximin* payoff level in the game. This means that (ii) for every player  $i \in I$ ,  $\wp_i(\mathbf{s}^x) \geq \wp_i(\mathbf{s}^{\vee i})$ . Combining (i) and (ii) gives

$$\min_{i \in I} [\wp_i(\mathbf{s}^y) - \wp_i(\mathbf{s}^{\vee i})] > \min_{i \in I} [\wp_i(\mathbf{s}^x) - \wp_i(\mathbf{s}^{\vee i})]$$

Hence,

$$\begin{aligned} \mathbf{s}^x \notin \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \} := \\ \{ \mathbf{s}^\tau \mid \forall \mathbf{s} \in \Sigma : \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \leq \min_{i \in I} [\wp_i(\mathbf{s}^\tau) - \wp_i(\mathbf{s}^{\vee i})] \} \end{aligned}$$

and so  $\mathbf{s}^x \notin \Sigma^\tau$ .

**Axiom 2 (Preservation of personal security):** *Team play cannot leave any player worse-off than his or her personal maximin payoff level in a game.*

Let  $s_i^\vee \in S_i$  be player  $i$ 's *maximin* (pure) strategy in a game. Let  $\mathbf{s}^\vee = (s_1^\vee, \dots, s_m^\vee)$  be a strategy profile where each player  $i \in I$  plays his or her maximin (pure) strategy  $s_i^\vee \in S_i$ . Every player preferentially ranks all strategy profiles that are associated with his or her *maximin strategy* at least as high as the strategy profile(s) associated with his or her *maximin payoff level* in the game. As such, for every player  $i \in I$ ,  $\wp_i(\mathbf{s}^\vee) \geq \wp_i(\mathbf{s}^{\vee i})$  (recall that  $\mathbf{s}^{\vee i}$  is a strategy profile that is associated with player  $i$ 's *maximin* payoff level in the game). Since the strategy profile  $\mathbf{s}^\vee \in \Sigma$  exists in every game, it follows that

$$\max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \} \geq 0$$

Therefore, if a strategy profile  $\mathbf{s}^x$  is such that, for some player  $i \in I$ ,  $\wp_i(\mathbf{s}^x) < \wp_i(\mathbf{s}^{\vee i})$ , then

$$\begin{aligned} \mathbf{s}^x \notin \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \} := \\ \{ \mathbf{s}^\tau \mid \forall \mathbf{s} \in \Sigma : \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \leq \min_{i \in I} [\wp_i(\mathbf{s}^\tau) - \wp_i(\mathbf{s}^{\vee i})] \} \end{aligned}$$

and so  $\mathbf{s}^x \notin \Sigma^\tau$ .

**Axiom 3 (Maximal mutual advantage):** *An outcome selected by a team has to be maximally mutually advantageous.*

The extent of individual advantage provided by a strategy profile  $\mathbf{s} \in \Sigma$  to a particular player  $i \in I$  is given by  $\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})$ . Let  $d(\mathbf{s})$  denote the extent of mutual advantage provided by the strategy profile  $\mathbf{s} \in \Sigma$  to the interacting players. Since the extent of mutual advantage is given by the number of units the strategy profile  $\mathbf{s} \in \Sigma$  advances all players' personal interests in parallel relative to the strategy profile(s) associated with each player's personal *maximin* payoff level in the game,

$$d(\mathbf{s}) = \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})]$$

and it is the case that, for every player  $i \in I$ ,  $\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i}) \geq d(\mathbf{s})$ .

Suppose that some strategy profile  $\mathbf{s}^x \in \Sigma$  is selected by a team that is not maximally mutually advantageous. In other words, there exists another strategy profile  $\mathbf{s}^y \in \Sigma$  such that  $d(\mathbf{s}^y) > d(\mathbf{s}^x)$ . It then follows from the above that

$$\min_{i \in I} [\wp_i(\mathbf{s}^y) - \wp_i(\mathbf{s}^{\vee i})] > \min_{i \in I} [\wp_i(\mathbf{s}^x) - \wp_i(\mathbf{s}^{\vee i})]$$

Hence

$$\begin{aligned} \mathbf{s}^x \notin \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \} := \\ \{ \mathbf{s}^\tau \mid \forall \mathbf{s} \in \Sigma : \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \leq \min_{i \in I} [\wp_i(\mathbf{s}^\tau) - \wp_i(\mathbf{s}^{\vee i})] \} \end{aligned}$$

and so  $\mathbf{s}^x \notin \Sigma^\tau$ .

**Axiom 4 (Indifference between equal extent of mutual advantage):** *If some outcome  $x$  is in a team's choice set and some other outcome  $y$  provides the same extent of mutual advantage as  $x$ , then  $y$  is in the team's choice set as well.*

Suppose that some strategy profile  $\mathbf{s}^x \in \Sigma$  is selected by a team and that some other strategy profile  $\mathbf{s}^y \in \Sigma$  provides the same extent of mutual advantage to the interacting players as the strategy profile  $\mathbf{s}^x$ . This means that  $d(\mathbf{s}^y) = d(\mathbf{s}^x)$  and that

$$\min_{i \in I} [\wp_i(\mathbf{s}^y) - \wp_i(\mathbf{s}^{\vee i})] = \min_{i \in I} [\wp_i(\mathbf{s}^x) - \wp_i(\mathbf{s}^{\vee i})]$$

Since  $\mathbf{s}^x$  is selected by the team,

$$\begin{aligned} \mathbf{s}^x \in \Sigma^\tau \equiv \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \} := \\ \{ \mathbf{s}^\tau \mid \forall \mathbf{s} \in \Sigma : \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \leq \min_{i \in I} [\wp_i(\mathbf{s}^\tau) - \wp_i(\mathbf{s}^{\vee i})] \} \end{aligned}$$

From above it follows that  $\mathbf{s}^y \in \Sigma^\tau$  as well.

## Appendix B

We here show that the presented rank-based function of team's interests is unique in satisfying the Axioms 1-4 introduced above. As in Appendix A, the discussion below is tailored to the case where the extent of individual and mutual advantage provided by different outcomes to the interacting players in games is measured using preferential rank values discussed on page 13. For a version of this proof where the extent of both types of advantage is measured using normalized payoffs discussed on page 16, all instances of  $\wp_i(\cdot)$  need to be replaced with  $u_i^*(\cdot)$ .

Suppose that some function  $\mathcal{G}^\tau(\Sigma) = \Sigma^\tau$  selects a subset from the set of all strategy profiles of the game  $\Gamma$  in a way that satisfies the Axioms 1-4 introduced above. From the definitions of individual and mutual advantage it follows that the extent of mutual advantage provided to the interacting players by some strategy profile  $\mathbf{s} \in \Sigma$  in a game is the least extent of individual advantage provided by that strategy profile across players and is given by

$$d(\mathbf{s}) = \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})]$$

(see the first paragraph of proof for Axiom 3 in Appendix A and the notation introduced in section 4.1). By Axiom 3, any strategy profile  $\mathbf{s}^\tau \in \Sigma^\tau$  selected by a team has to be maximally mutually advantageous, which means that the following condition must hold:

$$\forall \mathbf{s} \in \Sigma : \min_{i \in I} [\wp_i(\mathbf{s}^\tau) - \wp_i(\mathbf{s}^{\vee i})] \geq \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})]$$

This can be rewritten as

$$\begin{aligned} \mathbf{s}^\tau \in \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \} := \\ \{ \mathbf{s}^\tau \mid \forall \mathbf{s} \in \Sigma : \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \leq \min_{i \in I} [\wp_i(\mathbf{s}^\tau) - \wp_i(\mathbf{s}^{\vee i})] \} \end{aligned}$$

Suppose that some strategy profile  $\mathbf{s}^x \in \Sigma$  is selected by a team:  $\mathbf{s}^x \in \Sigma^\tau$ . From above,

$$\mathbf{s}^x \in \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \}$$

For any other strategy profile  $\mathbf{s}^y \in \Sigma$  that provides the same extent of mutual advantage to the interacting players as the strategy profile  $\mathbf{s}^x$ ,  $d(\mathbf{s}^y) = d(\mathbf{s}^x)$ , or, equivalently,

$$\min_{i \in I} [\wp_i(\mathbf{s}^y) - \wp_i(\mathbf{s}^{\vee i})] = \min_{i \in I} [\wp_i(\mathbf{s}^x) - \wp_i(\mathbf{s}^{\vee i})]$$

From this combined with the above it follows that

$$\mathbf{s}^y \in \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \}$$

Conversely, for any other strategy profile  $\mathbf{s}^z \in \Sigma$  that belongs to the set

$$\arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \}$$

it is the case that

$$\min_{i \in I} [\wp_i(\mathbf{s}^z) - \wp_i(\mathbf{s}^{\vee i})] = \min_{i \in I} [\wp_i(\mathbf{s}^x) - \wp_i(\mathbf{s}^{\vee i})]$$

which means that the strategy profile  $\mathbf{s}^z$  provides the same extent of mutual advantage to the interacting players as the strategy profile  $\mathbf{s}^x$ .

By Axiom 4, any strategy profile that provides the same extent of mutual advantage to the interacting players as the strategy profile  $\mathbf{s}^x \in \Sigma^\tau$ , or, equivalently, any strategy profile that belongs to the set

$$\arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \}$$

is in the team's choice set and, hence,

$$\Sigma^\tau \equiv \arg \max_{\mathbf{s} \in \Sigma} \{ \min_{i \in I} [\wp_i(\mathbf{s}) - \wp_i(\mathbf{s}^{\vee i})] \}$$

and  $\mathcal{G}^\tau \equiv \mathcal{F}^\tau$ .

## References

- Bacharach, M. [1999]: ‘Interactive Team Reasoning: A Contribution to the Theory of Co-operation’, *Research in Economics*, **53**, pp. 117–147.
- Bacharach, M. [2006]: *Beyond Individual Choice: Teams and Frames in Game Theory*, Princeton: Princeton University Press.
- Bardsley, N., Mehta, J., Starmer, C. and Sugden, R. [2010]: ‘Explaining Focal Points: Cognitive Hierarchy Theory versus Team Reasoning’, *The Economic Journal*, **120**, pp. 40–79.
- Battalio, R., Samuelson, L. and Van Huyck, J. [2001]: ‘Optimization Incentives and Coordination Failure in Laboratory Stag Hunt Games’, *Econometrica*, **69**, pp. 749–764.
- Colman, A. M., Pulford, B. D. and Lawrence, C. L. [2014]: ‘Explaining Strategic Coordination: Cognitive Hierarchy Theory, Strong Stackelberg Reasoning, and Team Reasoning’, *Decision*, **1**, pp. 35–58.
- Colman, A. M., Pulford, B. D. and Rose, J. [2008]: ‘Collective Rationality in Interactive Decisions: Evidence for Team Reasoning’, *Acta Psychologica*, **128**, pp. 387–397.
- Crawford, V. P., Gneezy, U. and Rottenstreich, Y. [2008]: ‘The Power of Focal Points Is Limited: Even Minute Payoff Asymmetry May Yield Large Coordination Failures’, *The American Economic Review*, **98**, pp. 1443–1458.
- Gauthier, D. [2013]: ‘Twenty-Five On’, *Ethics*, **123**, pp. 601–624.
- Gold, N. [2012]: ‘Team Reasoning, Framing and Cooperation’, in S. Okasha and K. Binmore (eds), ‘Evolution and Rationality: Decisions, Co-operation and Strategic Behaviour’, Cambridge: Cambridge University Press, chap. 9, pp. 185–212.
- Gold, N. and Sugden, R. [2007a]: ‘Collective Intentions and Team Agency’, *Journal of Philosophy*, **104**, pp. 109–137.
- Gold, N. and Sugden, R. [2007b]: ‘Theories of Team Agency’, in F. Peter and H. B. Schmid (eds), ‘Rationality and Commitment’, Oxford: Oxford University Press, chap. III.12, pp. 280–312.
- Kalai, E. [1977]: ‘Proportional Solutions to Bargaining Situations: Interpersonal Utility Comparisons’, *Econometrica*, **45**, pp. 1623–1630.
- Kalai, E. and Smorodinsky, M. [1975]: ‘Other Solutions to Nash’s Bargaining Problem’, *Econometrica*, **43**, pp. 513–518.
- Ledyard, J. O. [1995]: ‘Public Goods: A Survey of Experimental Research’, in J. H. Kagel and A. E. Roth (eds), ‘The Handbook of Experimental Economics’, Princeton: Princeton University Press, chap. 2, pp. 111–194.

- Luce, R. D. and Raiffa, H. [1957]: *Games and Decisions: Introduction and Critical Survey*, New York: Dover Publications, Inc.
- Misyak, J. and Chater, N. [2014]: ‘Virtual Bargaining: A Theory of Social Decision-Making’, *Philosophical Transactions of the Royal Society B*, **369**, pp. 1–9.
- Myerson, R. B. [1977]: ‘Two-Person Bargaining Problems and Comparable Utility’, *Econometrica*, **45**, pp. 1631–1637.
- Nozick, R. [1974]: *Anarchy, State, and Utopia*, New York: Basic Books, Inc.
- Nydegger, R. and Owen, G. [1974]: ‘Two-Person Bargaining: An Experimental Test of the Nash Axioms’, *International Journal of Game Theory*, **3**, pp. 239–249.
- Rankin, F. W., Van Huyck, J. B. and Battalio, R. C. [2000]: ‘Strategic Similarity and Emergent Conventions: Evidence from Similar Stag Hunt Games’, *Games and Economic Behavior*, **32**, pp. 315–337.
- Rawls, J. [1971]: *A Theory of Justice*, Cambridge, MA: Harvard University Press.
- Smerilli, A. [2012]: ‘We-Thinking and Vacillation Between Frames: Filling a Gap in Bacharach’s Theory’, *Theory and Decision*, **73**, pp. 539–560.
- Sugden, R. [1993]: ‘Thinking as a Team: Towards an Explanation of Nonselfish Behavior’, *Social Philosophy and Policy*, **10**, pp. 69–89.
- Sugden, R. [2000]: ‘Team Preferences’, *Economics and Philosophy*, **16**, pp. 175–204.
- Sugden, R. [2003]: ‘The Logic of Team Reasoning’, *Philosophical Explorations: An International Journal for the Philosophy of Mind and Action*, **6**, pp. 165–181.
- Sugden, R. [2011]: ‘Mutual Advantage, Conventions and Team Reasoning’, *International Review of Economics*, **58**, pp. 9–20.
- Sugden, R. [2015]: ‘Team Reasoning and Intentional Cooperation for Mutual Benefit’, *Journal of Social Ontology*, **1**, pp. 143–166.