# Superagency: Beyond an Individualistic Game Theory

Michael Bacharach,
University of Oxford

## Extended Abstract

In this lecture I suggest that some fundamental problems in game theory may be resolved by allowing that individual players can, without communication, coalesce to act like new, supraindividual players.

The study of interactive decision making has been pursued on two separate levels: the social level of social choice theory, cooperative game theory and DAI; and the individual level of noncooperative game theory. I predict that the two levels will stay apart not much longer. Game theory is mostly concerned with two fundamental features of rational interaction: isolation, and conflict of interest. DAI is mainly concerned with systems in which there is communication among agents with a common goal. As conflict and communication both dwindle, the two subjects converge on the same problem, the Silent Coordination Problem. At this interface between the subjects we find ouselves in a bare landscape dominated by a single large question: what is required to unite the rationality of the individual and the rationality of the group? This momentous question has been the object of surprisingly little attention, perhaps because this place has been so little visited.

What I shall call Schelling's Puzzle was how to explain what game theory could not, people's remarkable ability to coordinate, in pure coordination games such as Rendezvous, without a word passing between them [Sche60]. For many years Schelling's suggestion, that they did it by trying to 'focus' on a common point in the outcome space, and that this would lead them both to choose the most salient outcome, though unofficially accepted, remained unexplicated. A few years ago voices began to be heard urging that a theory of a game should include the players' own representations of their decision problem [Bach93, BaBe97, Rubi91]. One development in this spirit was the variable frame theory of games. This theory includes in the model of a game the bundle of classifications deployed by each player, her 'frame', assumed stochastic. Applying the theory to coordination games showed that, under a plausible general empirical condition on frames, each player's most probable representation of the game has the structure of the game of Hi-Lo:

|   | A | B |
|---|---|---|
| A | 2,2 | 0,0 |
| B | 0,0 | 1,1 |

Since the payoff-dominant equilibrium of this representation, (A, A), is also the outcome that comes out most salient under the stochastic framing process, this theory gives a game-theoretic answer to Schelling's puzzle.

But it left in its wake another, deeper, puzzle. It is obviously rational to play A in Hi-Lo; yet the proposition that A should be chosen is not derivable in classical game

theory. Everything is rationalizable, and (B, B) is an equilibrium. Call this the Hi-Lo Paradox. It's no good appealing to the salience of (2, 2), or to the payoff dominance of (A, A), or to equiprobability: these accounts of the rationality of A either merely reformulate the game as another case of Hi-Lo, or assume the principle that is to be shown, or treat one player as less than rational. The fundamental difficulty remains. It is that the social rationality of the Pareto principle is not reducible to the rationality of individual utility maximization.

It might be said: however engaging the Hi-Lo Paradox may be, Hi-Lo is only a limiting case – the fact that the discipline of Distributed AI has been almost wholly concerned with positive communication points to the empirical marginality of the Silent Coordination Problem. And it may seem that we have only to introduce a thread of communication into Hi-Lo, or indeed into any coordination problem, to get efficiency. But these appearances are deceptive: as long as we stick to the ground rules of classical game theory, communication is not sufficient to solve the coordination problem. Even if players with perfectly coincident interests are able to send messages to each other, the paradox survives intact, because truthtelling-and-accepting is merely the payoff-dominant equilibrium of another Hi-Lo game, in which falsetelling-and-rejecting is another.

But now for the good news: communication may not be sufficient for the resolution of Hi-Lo but, I argue, essentially it is not necessary either. Provided – as I shall for the moment suppose – computational costs are low enough, the role of communication can be substituted by a certain decision procedure in which each agent uses only her own initial information. I shall call this procedure *team reasoning*. In its most rudimentary form, *simple team reasoning*, each agent of a set (the 'team') computes the optimal profile of actions for the set, then chooses as her action her component in this profile. For some purposes it is useful to think of the simple team reasoner as simulating the computations and instructions of a 'phantom manager' in order to determine the answer to the question 'What would he have me do?' It has been suggested that footballers who improvise coordinated movements in the course of play may employ simple team reasoning. In Hi-Lo if both players simply team reason then both choose A. Sugden argued [Sugd93] that the optimality of (A, A) justifies team reasoning by each agent if it is common knowledge between them that both simply team reason.

If only some simply team reason, the group may not do well. The best move if all play their parts might risk conceding a goal if not all play their parts, Stag-Hunt-wise. Team reasoning must therefore take a different form when not all members of a group can be relied on to use it. Let $p(w)$ be the profile $p$ that is best for the group given that with probability $w$ each member will act in $p$ and with probability $1 - w$ will follow a default decision procedure. *Circumspect team reasoning* consists in computing $p(w)$, then choosing as your action your component in it. There is no need to require certainty that all are using team reasoning to assure its success. It is necessarily true – however unreliable co-members may be – that if everyone on team uses circumspect team reasoning then the profile it tells them to act in has an expected common payoff higher than (or equal to) tha of any other profile.

Later I shall say a little about the question of whether and if so when people team reason. But my primary concern is not with this empirical question, but with the normative question of the feasibility and efficacy of alternative decision procedures in

games. Because the informational inputs required by team reasoning and by classical, 'best reply' reasoning are the same (and the computational requirements likely to be lower), it is *open to* agents to team reason, just as it is open to them to use 'best reply' reasoning, and in some important cases we have already seen that the former is more efficacious.

The notion of circumspect team reasoning is easily generalized to the case in which each agent belongs to an arbitrary number of groups for which she may team reason: Bach99 studies the equilibria of the resulting 'unreliable team interactions' (uti's). The objectives of the different teams can be assumed to be anything, reflecting the variety of goals of the different groups and institutions to which, in real life, people are more or less strongly and constantly affiliated. An important special case, the *elementary uti*, is that in which each agent team reasons either for the whole group or for the singleton team consisting of herself alone.

It is shown in this theory that an interaction between $n$ 'unreliable teams' is behaviourally equivalent in a defined sense to a game between $n$ phantom managers, one per team. In particular, an elementary uti is behaviourally equivalent to an $n + 1$-player game. This theorem gives a precise game-theoretical sense in which team reasoning involves *sui generis* decision-makers.

Team reasoning is not at all the same thing as being motivated by the objective of the group. This is clear from the example of Hi-Lo. Take the group preference ordering to be the common individual ordering. Although, here, the individual player is motivated by the group objective, classic, individualistic reasoning does not, as we have seen, mandate A; while team reasoning does. In more general settings, a group goal and individual goals might diverge. In this case, team reasoning for the group goal involves not one but two transformations of the mental activitiy of a subject of classic game theory. The first is a payoff transformation, and the second the adoption of the decision procedure team reasoning. I shall return in a while to discuss circumstances in which both of these might occur.

First I want to chip away some more of the brick wall that has stood for so long between the domains of individual and social reasoning. Team reasoning has cons as well as pros: it enables coordination on a Pareto-efficient outcome; but on the other hand it involves duplication of computations, since each agent computes the entire profile but discards $n-1$ elements of it. These redundant computations may be extremely costly. It resembles the DAI system of Gene86 in which each deduces others' plans from the shared initial data, and it suffers from same limitations. A more thoroughgoing reform of the decision procedure of a group of agents who need to coordinate would feature communication and authority; it would go beyond team reasoning of the above kind to 'organizational reasoning'. By an organization I mean, roughly, a multi-agent system in which the agents are human agents with common objectives. As in Wern88, some agents compute plans or partial plans and send directives to others, whose decisions are governed by them. Forms of organization go from team reasoning through the master-slave relationship, plan passing, mutual cooperation, to societal cooperation. Who gets to send directives depends both on who observes different aspects of the environment, as in the 'theory of teams' [MaRa72] and on the need to avoid duplication of costly computations.

Just as it is open to players to team reason, it is open to them to reason as the members of an organization in which planning is distributed and relayed by directive. Consider the example of push-starting a car. It can happen that, without any discussion, the two pushers run behind and push, the driver begins to let the clutch in when the speed is high enough, and when the engine catches she makes a Gricean signal and the pushers stop. More generally, the possibility of 'spontaneous organization' follows at once from the fact that sending and acting on linguistic messages are, game theoretically, choices of strategies like any others. Hence adopting the particular roles of computation, message transmission, and directive execution specified in the optimal organizational design is, in cases of spontaneous organization, an application of team reasoning rather than an extension of it.

I conjecture that the opportunities the world presents for organization that both is complex and is produced spontaneously, by team reasoning, are limited, by bounds on human information processing capacities. Complex organization generally has a different aetiology, such as conscious design or incremental evolution.

Team reasoning forms a bridge between individual and group rationality. It happens in individual heads, but it harnesses the group's causal powers. It makes of the $n$ individual agents a single superagency. It is related to notions of collective intention discussed by Gilb88, Sear96 and others: at the end of a Sugden-style episode in which $n$ agents team reason to their common knowledge, their decisions together with their reasons for them place them in a state of common intention in Bratman's sense (though not in Gilbert's since they make no explicit commitments to each other).

Where group and individual interests are identical, team reasoning is a means of achieving Pareto optimality without prior agreement or antecedent organization. In other cases it is a means by which groups' goals differing from individuals' goals may be achieved if and when the individuals engage in it. This they might do under duress, or because they are paid to, or because they are socialized to, but it can also happen through a voluntary process, 'group identification'. It is plausible that in addition to its well-documented effects such as pro attitudes towards in-group members, group identification primes team reasoning. There is evidence that it is associated with a 'frame' characterized by 'we' concepts, so that it is likely that a group identifier asks herself not 'What should I do?' but 'What should we do?'. This question calls for the computation of the optimal profile for the group, which is the first step of team reasoning. The second step, adopting the questioner's component in the optimal profile as her decision, is then a simple application of the deontic rule of inference $(A \rightarrow B) \vdash (\Box A \rightarrow \Box B)$.

## REFERENCES

Bacharach, Michael (1993), 'Variable Universe Games', in Binmore, Ken *et al* (eds), Frontiers of Game Theory (MIT Press)

Bacharach, Michael (1999), 'Interactive team reasoning: A contribution to the theory of cooperation', Research in Economics 53, 117-47

Bacharach, Michael and Bernasconi, Michele (1997), 'The Variable Frame Theory of Focal Points: An Experimental Study', Games and Economic Behavior 19, 1-45

Gesenereth, Michael *et al* (1986), 'Cooperation without Communication', Proceedings AAAI-86

Gilbert, Margaret (1988), On Social Facts (Routledge)

Marschak, J. and Radner, R. (1972), Economic Theory of Teams (Yale University Press)

Rubinstein, Ariel (1991), 'Comments on the Interpretation of Game Theory', Econometrica 59, 909-24

Schelling, Thomas C. (1960), The Strategy of Conflict (Harvard University Press)

Searle, John (1996), The Construction of Social Reality

Sugden, Robert (1993), 'Thinking as a Team: Towards an Explanation of Nonselfish Behavior', Social Philosophy and Policy 10, 69-89

Werner, Eric (1988) 'Toward a theory of communication and cooperation for multiagent planning', in Moshe Vardi (ed), Proceedings of the Second Conference on Theoretical Aspects of Reasoning about Knowledge (Morgan Kaufmann)