"Beliefs, Intentions and Emotions:
Old versus New Psychological Game Theory"

by

Jeffrey Paul Carpenter
Peter Hans Matthews

January, 2003

# Beliefs, Intentions and Emotions:
## Old versus New Psychological Game Theory

(Invited comment on Andrew Colman, Cooperation, psychological game theory, and limitations of rationality in social interaction, forthcoming, *Behavioral and Brain Sciences*)

Jeffrey P. Carpenter
Department of Economics
Middlebury College
Middlebury, Vermont 05753

jpc@middlebury.edu

Peter Hans Matthews
Department of Economics
Middlebury College
Middlebury, Vermont 05753

peter.h.matthews@middlebury.edu

3 January 2003

Abstract: We compare Colman's proposed "psychological game theory" with the existing literature on psychological games (Geanakoplos, Pearce and Stachetti 1989), in which beliefs and intentions assume a prominent role. We also discuss the experimental evidence on intentions, with a particular emphasis on reciprocal behavior, as well as some recent efforts to show that such behavior is consistent with social evolution.

Andrew Colman's provocative paper is a manifesto of sorts, a call to build a new, *psychological*, game theory based on "nonstandard assumptions." Our immediate purpose is to remind readers that the earlier work of Geanakoplos, Pearce and Stachetti (1989) or GPS, which the paper cites but does not discuss in much detail, established the foundations for a theory of "psychological games" that achieves at least some of the same ends. Our brief review of GPS and some of its most important descendants – in particular, the work of Rabin (1993) and Falk and Fischbacher (2000) – will also allow us to elaborate on some of the connections between psychological games, experimental economics and social evolution.

The basic premise of GPS is that payoffs are sometimes a function of both actions *and* beliefs about these actions, where the latter assumes the form of a subjective probability measure over the product of strategy spaces. If these beliefs are "coherent" – that is, the information embodied in second order beliefs are consistent with the first order beliefs, and so on – and this coherence is common knowledge, then the influence of second (and third and fourth …) order beliefs can be reduced to a set of common first order beliefs. That is, in a two player psychological game, for example, the utilities of *A* and *B* are functions of the (perhaps mixed) strategies of each and the beliefs of each about these strategies. A psychological Nash equilibrium or PNE is then a strategy profile in which given their beliefs, neither A nor B would prefer to deviate *and* these first order beliefs are correct. If these augmented utilities are continuous – not an innocuous restriction in this framework – then all normal form psychological games must have at least one PNE.

The introduction of beliefs provides a natural framework for modeling the role of *intentions* in strategic contests, and this could well prove the most important application of GPS. It is obvious that intentions matter to decision-makers – consider the legal difference between manslaughter and murder – and that game theorists would do well to heed the advice of Colman and others who advocate a more behavioral approach.

For a time, it was not clear whether or not the GPS framework was tractable. Rabin (1993), which Colman cites as an example of behavioral, rather than psychological, game theory, was perhaps the first to illustrate how a normal form psychological game (in the

sense of GPS) could be derived from a "material game" with the addition of parsimonious "kindness beliefs." In the standard two person prisoner's dilemma or PD, for example, he showed that the "all cooperate" and "all defect" outcomes could *both* be rationalized as PNEs.

Consistent with the experimental evidence, the structure of beliefs in Rabin (1993) underscores the importance of intentions. The "all defect" PNE in his transformed PD occurs when agents believe, for whatever reason, that other agent(s) are ill-intentioned. In contrast, for a population of well-intentioned agents with common knowledge of these intentions, where this knowledge could be the result of past interaction, the "all cooperate" PNE arises.

As Rabin [1993] himself notes, this transformation of the PD is *not* equivalent to the substitution of altruistic agents for self-interested (in the traditional sense) ones: the "all defect" outcome, in which each prisoner believes that the other(s) will defect, could not otherwise be an equilibrium. This is an important caveat to the recommendation that we endow economic actors with "nonstandard reasoning processes," and prompts the question: What observed behavior will the "new psychological game theory" explain that an old(er) GPS-inspired one cannot? Or, in narrower terms, what are the shortcomings of game theoretic models that incorporate the role of intentions, and therefore such emotions as surprise or resentfulness?

The answers are not obvious, not least because there are so few examples of the transformation of material games into plausible psychological ones, and almost all of these share Rabin's (1993) emphasis on kindness and reciprocal behavior. It does seem to us, however, that to the extent that Colman's "nonstandard reasoning" can be formalized in terms of intentions and beliefs, there are fewer differences between the old and new psychological game theories than first seems.

There is considerable experimental evidence that intentions matter. Consider, for example, the experiment described in Falk, Fehr and Fischbacher (2000), in which a first mover can either give money to, or take money away from, a second mover, and any money

given is tripled before it reaches the second mover, who must then decide whether to give money back, or take money from, the first mover. Their analysis suggests that there is a strong relationship between what the first and second movers do: in particular, the more the first mover gives (takes) the more the second mover takes gives (takes) back.

Falk *et al* (2000) find that first mover giving (taking) is interpreted as a friendly (unfriendly) act, and that these intentions matter. Absent the influence of beliefs or intentions on utilities, there would be a single Nash equilibrium in which the first mover takes as much as possible from the second because she "knows" that the second has no material incentive to retaliate. While this behavior can also be supported as a PNE, so can that in which the first mover gives and expects a return and the second mover understands this intention and reciprocate. When the experiment is changed so that the first mover's choice is determined randomly, so that there are no intentions for the second mover to impute, the correlation between first and second mover actions collapses. We see this as (a) further evidence that beliefs – in particular, intentions – matter to decision makers but also (b) that once these beliefs have been incorporated into strategic models, a modified "rational choice framework" is still useful.

The old psychological game theory can also accommodate at least some of Colman's legitimate concerns about rational choice in sequential games. Geanakoplos (1996), for example, draws on GPS to propose new solutions, or at least shed new light on, two famous problems, the Hangman's Paradox and Newcomb's Paradox. Both exploit the fact that the standard backward induction arguments become inappropriate because standard extensive form nodes are an incomplete description of the state of the world where utilities are also a function of beliefs/intentions.

Building on both GPS and Rabin (1993), Dufwenberg and Kirchsteiger (1998) allow for the revision of beliefs, and therefore perceptions of kindness, in extensive form games, which allows them to formalize the notion of "sequential reciprocity."

In a similar vein, Falk and Fischbacher (2000) derive a variation of Rabin's (1993) "fairness equilibrium" for a number of simple but important extensive form games, with

results that are also consistent in broad terms with the experimental evidence. The simplest of these is perhaps the ultimatum game, in which a first mover offers some share of a pie of known size to a second mover who must then accept or reject the proposal. With kindness functions similar to Rabin's (1993), Falk and Fischbacher (2000) show that the ultimatum game has a unique PNE that varies with the "reciprocity parameters" of proposer and responder. Furthermore, this equilibrium is consistent with the observations that the modal offer is half the surplus, that offers near the mode are seldom rejected, and that there are few low offers that are consistent with the subgame perfect equilibrium, and that most of these low offers are rejected.

This result does *not* tell us, though, whether this outcome is consistent with the development of reciprocal intentions or norms over time or, in other words, whether social evolution favors those with "good intentions." To be more concrete, suppose that the proposers and responders in the ultimatum game are drawn from two distinct populations and matched at random each period, and that these populations are (to start, at least) heterogeneous with respect to intention. Could these intentions survive some sort of "selection mechanism" based on differences in *material* outcomes? Or do these intentions impose substantial costs on those who have them?

There are still no definitive answers to these questions, but the results in Binmore, Gale and Samuelson (1995) or BGS hint that such intentions will sometimes be robust. BGS consider a "miniature ultimatum game" with two offers (fair and selfish) and two responses (accept and reject) and assume that the shares of proposers and responders who are committed to these pure strategies evolve on the basis of the standard replicator dynamic. There are *two* stable equilibria within this framework: the first corresponds to the subgame perfect equilibrium – all proposers are selfish, and all responders accept these selfish offers – but in the second, all proposers are fair and a substantial but indeterminate share of responders would turn down an unfair offer. Furthermore, these dynamics can be rationalized as a form of social or cultural learning: BGS the role of aspiration, but evolution toward fair outcomes is also consistent with imitation (Björnerstedt and Weibull 1996). It is tempting, then, to interpret the second BGS outcome as a Falk and Fischbacher (2000) "fairness equilibrium."

This interpretation is somewhat premature, however. On one hand, it isn't clear that social evolution can select beliefs or intentions: Gintis (2000) asserts, for example, that the subjective beliefs characteristic of all GPS-based models "lack explanatory power" in evolutionary models of reciprocal (or other) behavior. It is therefore important to note that Gintis (2000) is also able to reproduce the Falk and Fischbacher (2000) result when each agent holds a well-defined fairness norm. On the other hand, because BGS limit attention to pure strategies in the material game, the role of intentions or for that matter norms is at best implicit. Indeed, as the title of the BGS paper reminds us, one *could* characterize the behavior on the path(s) that lead to fair outcomes as "learning to be imperfect" rather than a consequence of the fitness-enhancing predilection for pro-social norms or beliefs.

All of this said, we share most of Colman's concerns with standard game theoretic arguments, and suspect that psychological game theorists, both old *and* new, will have much to contribute to the literature in the decades ahead.

Binmore, K., Gale, J. & Samuelson, L. (1995) Learning to be imperfect: the ultimatum game. *Games and Economic Behavior* 8: 56-90.

Björnerstedt, J. and Weibull, J. (1996) Nash equilibrium and evolution by imitation, in K. Arrow, E. Colombatto, M. Perlman & E. Schmidt, eds, *The Rational Foundations of Economic Behavior*. New York: Macmillan.

Dufwenberg, M. & Kirchsteiger, G. (1998) A theory of sequential reciprocity. Tilburg CentER for Economic Research Discussion Paper No. 9837.

Falk, A. & Fischbacher, U. (2000) A theory of reciprocity. Institute for Empirical Research in Economics Working Paper No. 6.

Falk, A., Fehr, E. & Fischbacher, U. (2000) Testing Theories of Fairness – Intentions Matter. Institute for Empirical Research in Economics Working Paper No. 63.

Geanakoplos, J. D. (1996) The hangman's paradox and Newcomb's paradox as psychological games. Cowles Foundation Discussion Paper No. 1128.

Geanakoplos, J. D., Pearce, D. & Stachetti, E. (1989) Psychological games and sequential rationality. *Games and Economic Behavior* 1: 60-79.

Gintis, H. (2000) *Game Theory Evolving*. Princeton: Princeton University Press.

Rabin, M. (1993) Incorporating fairness into game theory and economics. *American Economic Review* 83: 1281-1302.