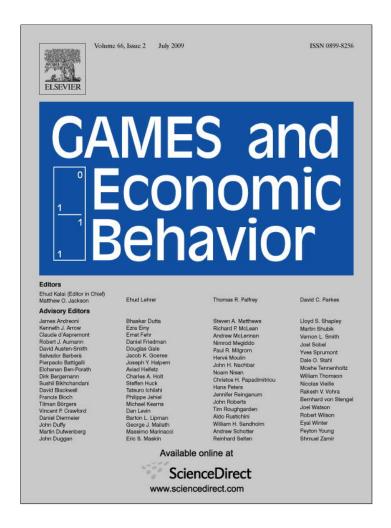
Provided for non-commercial research and education use. Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

http://www.elsevier.com/copyright



Available online at www.sciencedirect.com



Games and Economic Behavior 66 (2009) 761-774



www.elsevier.com/locate/geb

# Endogenous games and equilibrium adoption of social norms and ethical constraints \*

John P. Conley a, William Neilson b,\*

<sup>a</sup> Department of Economics, Vanderbilt University, Nashville, TN 37235, USA

<sup>b</sup> Department of Economics, University of Tennessee, Knoxville, TN 37996, USA

Received 7 March 2007

Available online 7 October 2008

#### **Abstract**

We consider a situation in which games are formed endogenously in two senses: (1) there is a pregame in which agents choose to learn a subset of all feasible strategies and can then employ only these strategies in subsequent play, and (2) agents choose their game partners through a costly search process. We show that at any subgame perfect equilibrium, agents will constrain their action sets in the pregame in such a way that a single social norm prevails. Thus, all agents in a society will abide by the same ethical standard, although what standard this will be cannot be predicted. We also show that these are essentially the only SPE outcomes. We suggest that this provides at least a partial explanation for experimental observations that agents apparently choose strategies that do not maximize their payoffs.

© 2008 Elsevier Inc. All rights reserved.

JEL classification: C72; C78; Z13

Keywords: Behavioral economics; Endogenous games; Bilateral bargaining; Prisoners' dilemma; Social norms

## 1. Introduction

Experiments and everyday observations show that in many situations, people do not seem to behave as classical game theory predicts. Many explanations have been offered in the literature. If people are boundedly rational, for example, it may be difficult for them to understand their best interests in complicated situations. While certainly true in the abstract, this does not help us understand why agents continue to follow apparently suboptimal strategies even in very simple strategic situations such as ultimatum or prisoners' dilemma games. Alternatively, it may be that agents are basing their actions on rules of thumb, a concern for the welfare of others, or on a prevailing social norm as behavioral economists suggest.

<sup>\*</sup> We wish to thank Ehud Kalai, Herve Moulin, Simon Wilkie, and the participants of PET06, Hanoi for their comments and suggestions. We also thank an anonymous referee who pointed out a major simplification for the type of dominance solvable stage games considered in this paper. We take responsibilities for all remaining errors.

<sup>\*</sup> Corresponding author at: University of Tennessee, Department of Economics, 523 Stokely Management Center, Knoxville, TN, United States. *E-mail addresses:* j.p.conley@vanderbilt.edu (J.P. Conley), wneilson@utk.edu (W. Neilson).

It would be foolish to deny these are important elements affecting human behavior. For example, professionals, even professors, devote considerable time and energy to meet self-imposed standards even when there is no expectation of reward. We do not expect our doctor to lie to us about our condition in order to run up the bill despite our inability to verify or even understand his diagnosis, for example. We count on his sense of professionalism or ethics. Why don't these agents follow their apparent self-interest? Whatever the reasons, the standard game theoretic approach clearly misses something in its description of human behavior.

In the real world, individuals can often choose the agents with whom they interact. For example, people choose their friends, spouses, business partners, and coauthors. Moreover, these choices are often informed by observable features of the members of the pool of potential partners. Sometimes these decisions are based solely on what one might view as utility-maximizing considerations, such as when an individual chooses a spouse based on attributes he or she would like to "consume," or when a researcher chooses a coauthor based on knowledge complementarities. Other times, however, these decisions are based on "strategic" considerations in the sense that pairing with one partner may entail playing a game in a more favorable way than pairing with another. For example, when choosing a coauthor, a researcher might prefer someone with a strong sense of personal responsibility. Such a partner, in effect, does not know how to free-ride and thus cannot play this strategy even though it is, in principle, a feasible choice.

The key idea is that when individuals can accept or reject matches, agents will find it in their interests to make themselves attractive partners. Thus, when partnership decisions are based on strategic considerations, agents should take steps to make themselves someone that others would like to play a game against.

To make this more concrete, consider the marriage market as an example. As we say above, marriage is a partner-ship which is entered into voluntarily. As in the Rubinstein and Wolinsky (1985) model, people meet many potential partners and incur search costs if they reject a match and seek an alternative partnership. While many marriages dissolve when one or the other partner feels that the relationship is no longer in their interests, we also have many examples of apparently selfless behavior. Spouses stick by one another even when one is injured, loses a job, becomes addicted to a substance, etc. Similarly, committed partners maintain fidelity even when there is little chance of a dalliance being discovered. Of course, the most immediate explanation for this is that these are ethical or moral choices. The question is: how can we understand this within the framework of rationally maximizing agents?

In searching for marriage partners, one of the things we try to find out about is the "character" of our potential mate: will this person be loyal, will he or she stick by us in difficult times, etc. Promises are not good enough since they are not subgame perfect. Reputation is of limited value as this game is not very frequently repeated. Partners therefore go to significant effort to signal they have credibly constrained their action sets to preclude these types of behavior. This might be through the vehicle of claiming affiliation with a particular religion, showing that one was raised by a family that disapproves of such behavior, demonstrating selfless loyalty to one's friends or community in some way, or even by showing a lack of social facility to play the field. In short, the notion we explore here is that certain actions are not available to you because your ethics forbid it (religious prohibitions against divorce) or you have never learned them (I simply have no idea how to attract the opposite sex).

In the context of this example, it is clear why agents would choose to constrain their strategy sets before seeking a mate. If a man, for example, wants to match high in the marriage market, he has to offer something of value in return. He should take great pains to signal that he is not the kind of guy who will cheat, seek a trophy wife in 10 years, or fail to work hard to provide for the children. By so doing, he can attract a wife with similar ideas who will continue to contribute to the marriage even if he loses his job, his rippling six-pack abs, or his health. As marriage is a positive sum game, they can divide a larger surplus by choosing to forgo the possibility of pursuing strategies that might be in their interests in some subgames.

The purpose of this paper is to explore more formally the idea of selecting partners for strategic reasons and the implications this has for the structure of the games that people play. To do this, we construct what we call an *endogenous game*. An endogenous game begins with a symmetric stage game, in our case a prisoner's dilemma, and a population of identical players. In the first period of the game, which we call a *pregame*, each player commits to a subset of the action space of the base game, and we call this chosen subset a player's *list*. In subsequent periods players are randomly matched, observe their matched opponents' lists, and then decide whether to play against that opponent or pay a waiting cost and get rematched next period. When a pair decides to play, they each choose single

<sup>&</sup>lt;sup>1</sup> Dekel et al. (2007) examine an evolutionary model in which agents with potentially different preferences are matched. When preferences are fully observable agents know their matched partners preferences, and when preferences are unobservable agents only know the population

actions from their lists, collect their payoffs, and leave the game. While the continuation game (after the pregame) is similar in spirit to the Rubinstein and Wolinsky (1985) two-sided bargaining model, the addition of the pregame makes the endogenous game considered here distinctive.

The game structure we explore is "endogenous" in two senses. Recall a standard game is defined by three elements: a set of players, an action set for each player, and a payoff function for each player. These are all taken as exogenous. In particular, agents are dropped into the game and must choose an action from the specified action set. In endogenous games, the player set arises though strategic choices of agents. No agent is forced to enter a stage game. In addition, the action set is chosen by agents in the pregame, and so this is endogenous as well. Thus, two out of three of the basic elements of a game arise endogenously in our setup.<sup>2</sup>

We begin by considering a prisoners' dilemma game. We show that if we restrict attention to pure strategies, only two robust outcomes are possible: all agents learn to cooperate (only) or all agents learn either only to defect or both to cooperate and defect but end up always choosing to defect.<sup>3</sup> If we modify this game to allow a finite set of partially cooperative behaviors then only two classes of outcomes are possible. When the set of partially cooperative strategies is coarse enough relative to the cost of search, then any outcome in which all agents learn the same, unique strategy is an equilibrium. More interestingly, these social norm outcomes are the only equilibria. It is not possible for agents to choose different lists or play different strategies in equilibrium. However, if the set of strategic alternatives becomes fine enough, then the game unravels in the sense that the only equilibrium is for all agents to learn and play the noncooperative strategy.

We also consider an "anti-prisoner's dilemma" in which "cooperation" is a dominant and Pareto optimal strategy. We now find that social norms that are strictly worse than the dominant strategy outcome for all agents can be supported as equilibria. Thus, when games are endogenous, we may see harmful social norms emerge. What drives agents to participate in this apparently self-defeating behavior is that no one will interact with an agent who fails to follow the prevailing social norm. In a sense, peer pressure supports undesirable social outcomes. The anti-prisoner's dilemma is also useful in making two other points about endogenous games clear. First, the pregame does not in any sense facilitate beneficial cooperation (as would a system of *ex ante* binding contracts, for example). Rather, the endogenous game structure is a neutral device that turns out to drive populations of agents to adopt a single social norm of behavior which may be better or worse for agents than the payoffs they would receive in the associated one-shot game. Second, endogenous games do not give a folk theorem like we see in repeated games. In particular, only symmetric outcomes are possible, and more importantly, equilibrium outcomes can be dominated by the minmax payoff.

The addition of a pregame allows for two new interpretations of how social norms arise. First, in our model, defections from the social norm result from agents adopting lists with strategies that are worse for the group than the social norm strategies. Agents then punish these defectors by refusing to play against them. In essence, defections from the social norm, if not advantageous to the group, result in ostracism of the defector.<sup>4</sup> This form of punishment contrasts with what we see in the repeated games literature where punishment takes place after a defection in the stage game, as in Coleman (1990), Kandori (1992), Ellison (1994), Neilson (1999), and Dal Bo (2007).<sup>5</sup>

A second interpretation arises from thinking about the choice of a list, and therefore the adoption of a social norm, as an investment in social capital. Because the investment comes at the beginning of the game, it is subject to holdup problems (e.g. Peters and Siow, 2002, and Cole et al., 2001). These holdup problems are mitigated by the ability to refuse to play against someone who underinvests in social capital. Similar phenomena arise in Burdett and Coles (2001) and Baker and Jacobsen (2007). Burdett and Coles consider a marriage market in which singles invest in order to make themselves more attractive, and find that bad social norm equilibria arise in which singles overinvest in attractiveness. Baker and Jacobsen examine a marriage matching problem that incorporates two elements, a holdup

distribution of preferences. Our paper has a different information structure, with agents observing their matched partners' lists but not the population distribution of lists. Rather, agents infer the population distribution of lists from the equilibrium strategies.

<sup>&</sup>lt;sup>2</sup> Dekel et al. (2007) explore implications of relaxing the third assumption, fixed payoff functions, using an evolutionary framework.

<sup>&</sup>lt;sup>3</sup> We will ignore for now an additional type of equilibrium that would not survive farsighted behavior by agents.

<sup>&</sup>lt;sup>4</sup> Hirshleifer and Rasmusen (1989) analyze a different model of ostracism in which current defections from the group-cooperative strategy are punished by future ostracism from the group.

<sup>&</sup>lt;sup>5</sup> Also see Axelrod (1986) who shows that social norms, such as cooperation in the prisoners dilemma, can arise out of evolutionary processes. Extensions of Axelrod's work include Bendor and Swistak (2001) and Sugden (2004).

<sup>&</sup>lt;sup>6</sup> Jansen (2004) explores similar phenomena in a job-matching setting. An overview of the search-matching literature can be found in Rogerson et al. (2005).

problem arising because of human capital investment coming before the marriage, and a cooperation problem caused by the joint valuation of human capital. They illustrate how custom can help mitigate the holdup problem.<sup>7</sup>

The paper proceeds as follows. In Section 2 we describe the model. In Sections 3 and 4 we give our results for the prisoner's dilemma and the anti-prisoner's dilemma, respectively. In Section 5 we discuss the existing literature as it relates to these results in more detail. We make a few remarks about our modeling approach in Section 6, and Section 7 concludes.

### 2. The model

Consider an economy with I agents indexed  $i \in \{1, ..., I\} \equiv \mathcal{I}$ . Each agent has a finite set of strategies  $\mathcal{X}$  available for him to learn. Agents receive payoffs by finding a partner with whom to play a bilateral game. The payoff to each agent is given by the function  $F : \mathcal{X} \times \mathcal{X} \mapsto \mathfrak{R}$ . Thus, the payoff to agent i choosing strategy  $x_i$  is  $F(x_i, x_i)$ .

We assume a two-part game. The first part we will call the *pregame*. Here, each agent chooses a subset of strategies to learn  $\ell_i \subseteq \mathcal{X}$ . We will refer to  $\ell_i$  as agent *i*'s *list*. Chosen lists are private information during the pregame, but are revealed on a player-by-player basis later in the game as described below. Denote the set of all possible lists as  $\mathcal{L}$  (that is, all possible nonempty subsets of  $\mathcal{X}$ ).

The second part of the game is a multistage matching game similar to the one described by Rubinstein and Wolinsky (1985). In each round of the stage game, agents are randomly matched with another agent, and each agent observes his partner's list. They then can choose one of two options: play or not play.

If at least one agent decides to Not play (N), they receive no payoff but pay a *delay cost* of z > 0 and search for a new partner in the next round of play. If both agents decide to Play (P), they choose strategies from their list, settle on a Nash equilibrium and retire from the game. Thus, the stage game really has two substages, one in which the agents decide to play or not, and the second substage where they decide on what strategy to employ if they do decide to play together. We will carry these stages in the proofs of our results, but in the interest of simplicity, will not create notation to describe this explicitly. Thus, choosing a strategy x implies that one chose P in first substage, and choosing a strategy N implies that one would play the dominant strategy in the second substage if it were reached.

We will also assume that when an agent retires, he is replaced by an agent with the same name who has chosen to play in the same way. While this is a strong assumption from a theoretical standpoint, it will become clear that in the applications we consider, it has very little bite.

Formally, an agent's strategy in the second part of the game is a mapping from the strategy list of the agent with whom he is matched to his own list plus N. Since matching is random and anonymous and the population of agent types is stable across periods, we assume that history does not influence strategic choice. Let agent *i*'s second stage strategy be denoted  $g_i : \mathcal{L} \mapsto \ell_i \cup N$ . Let  $\mathcal{G}$  denote the set of all such mappings. We will refer to  $g_i$  as agent *i*'s stage game strategy.

An agent's strategy therefore consists of a pair  $(\ell_i, g_i) \in (\mathcal{L}, \mathcal{G})$ . Collectively, we will denote this as  $s_i \in \mathcal{S}$ . A *strategy profile* for the game is denoted:  $S \equiv (s_1, \ldots, s_I) \in \mathcal{S} \times \cdots \times \mathcal{S}$ . It will be useful to refer to the profile of lists and stage game strategies separately on occasion. We denote these  $L = (\ell_1, \ldots, \ell_I)$  and  $G = (g_1, \ldots, g_I)$ . We will also use the notation  $s_{-i}$  to denote the strategy profile for all agents excluding agent i. Finally, let  $\bar{L}$  denote all the possible lists that are *not* chosen by any agent given the strategy profile S. Thus, no list  $\bar{\ell} \in \bar{L}$  should ever be seen by any agent in any stage game on the equilibrium path.

We will use the standard notion of Subgame Perfect Equilibrium (SPE) in the paper.

In the next two sections, we will maintain the two assumptions that population size, I, is even and  $I \ge 4$ . The first is because this is a matching game, and if the population were odd, we would have to include the possibility of not finding a potential match in any given period. This would needlessly complicate our model. The second is because if I = 2 agents would always be matched with the same player each period which would make the endogeneity of the game degenerate.

<sup>&</sup>lt;sup>7</sup> As with much of the search-matching literature, both of these papers analyze steady states. Consequently, we contribute to the search-matching literature on hold-up problems by allowing for strategic equilibria rather than steady-state equilibria.

<sup>&</sup>lt;sup>8</sup> We thank an anonymous referee for suggesting this.

#### 3. Endogenous play in a granular prisoners' dilemma game

| Prisoners' dilemma |           |        |  |
|--------------------|-----------|--------|--|
|                    | Cooperate | Defect |  |
| Cooperate          | 10, 10    | -3, 12 |  |
| Defect             | 12, -3    | 0, 0   |  |

In this section, we begin by defining a standard symmetric prisoners' dilemma game. The strategy set for each agent is  $\{c, d\}$  with the payoffs given by:

$$F^{pd}(c,c) = 10$$
,  $F^{pd}(d,c) = 12$ ,  $F^{pd}(c,d) = -3$ ,  $F^{pd}(d,d) = 0$ .

Note that the marriage game discussed in the introduction can be mapped onto this. Cooperating in this context means making contributions to the household (washing dishes, cooking dinner) that provide public benefits to both agents. Defecting means free riding on the contributions (if any) of one's spouse.

Rational players, of course, will never contribute to a marriage in one shot play. The main point of this section is to show that when games are formed endogenously, making contributions to the collective becomes rational. In addition, we show that the equilibrium contribution levels will be equal for both agents. Thus, endogenous games enforce the emergence of various different social norms of symmetric and reciprocal behavior. This contrasts with repeated versions of this game in which a folk theorem holds and asymmetric contributions might be seen in equilibrium.

We will restrict attention to pure strategy equilibria; however, we wish to allow agents richer strategy sets than simply cooperation and defection. We will therefore modify the PD game above to allow agents to choose from a finite set of evenly spaced partially cooperative strategies. There are a number of ways to do this, but the most direct is to imagine pure strategies that have payoffs equivalent to mixed PD strategies. Formally:

$$\mathcal{X} = \{0, \epsilon, 2\epsilon, \dots, (n-1)\epsilon, 1\}$$

where  $\epsilon = 1/n$ . The payoffs are:

$$F(x_i, x_j) = x_i x_j F^{pd}(c, c) + x_i (1 - x_j) F^{pd}(c, d) + (1 - x_i) x_j F^{pd}(d, c) + (1 - x_i) (1 - x_j) F^{pd}(d, d).$$

Note that we interpret these as partially cooperative pure strategies and just construct a game whose payoffs happen to correspond to those of a mixed strategy prisoners' dilemma. Since we don't literally contemplate mixed strategies, in practice these partially cooperative strategies might mean "learn how to wash dishes, but not how to do laundry" or "have a religious belief in the tithe but doubts about the immorality of birth control." Such agents would only be able to make bounded contributions to the collectives they joined. In any event, this is not material to the results. All we need is that there be a granularity in the strategy choices of the agents. We discuss below the outcome when agents have a continuum of strategies (either pure or mixed).

Of course, the analysis we carry out below applies immediately to a more general class of games than the granular prisoner's dilemma. Specifically, let the strategy set for each agent be  $\mathcal{X}=0,\ldots,k$ , and the payoff function satisfies F(i,j)>F(i',j) for all i'>i, and F(i,j)>F(i,j') for all j'< j. The first condition implies that an agent's strategies are ordered by dominance, with action 0 being the dominant strategy and action i dominating action i' whenever i'>i. The second condition states that switching from one strategy to a dominating one hurts a player's partner. This makes higher-numbered strategies more cooperative (or more generous) because they generate higher payoffs for the partner. A generalized form of the granular prisoner's dilemma results when F(i,i)>F(i',i') whenever i'< i. We will continue to state our results using the granular PD game above, however, since it is widely studied and makes the economic and behavioral implications of our results easy to talk about and interpret. We also note that it is this dominance solvability that allows us to use the standard notion of SPE. Without this, we would have to worry about agents' beliefs about how agents would play in the stage games, and how this in turn would affect the lists that were optimal to choose in the pregame. We are indebted to an anonymous referee for pointing out both of these observations, the latter of which allowed us to greatly simplify the paper.

**Theorem 1.** If z is sufficiently small, then for any  $\hat{x} \in \mathcal{X}$  there exists an SPE in which  $\ell_i = \{\hat{x}\}$  for all agents  $i \in \mathcal{I}$ .

**Proof.** We show this by backwards induction. Suppose two agents with these lists were matched. If they agree to play, then they each have only one possible strategy to play and so they get a payoff of  $f(\hat{x}, \hat{x})$ . Given that this is the best

they could do in any rematching as well, it is a best response to agree to play rather than wait to play with another agent in a future round and pay a delay cost. Thus all agents will agree to play when matched. Finally, we must show that no agent could do better by choosing a different list in the pregame.

We consider three cases. Suppose first that an agent deviated by choosing a list  $\ell_i$  such that  $\min \ell_i = \bar{x} < \hat{x}$  (that is, with a least cooperative element involving less cooperation that  $\hat{x}$ ). Since in any stage game, the only best response is for an agent to play his least cooperative strategy, and all other agents are identical and can only play  $\hat{x}$  in any stage game, agent i's best response is to agree to play whenever matched and play  $\bar{x}$ . However, if the cost of delay z is small, all other agents would be better off declining to play when matched with i, which would lead to payoff  $F(\hat{x}, \bar{x}) < F(\hat{x}, \hat{x})$ , in favor of waiting for a more cooperative player and receiving  $F(\hat{x}, \hat{x})$  in a future round. Thus, agent i will never find a match and will pay the delay cost z each round as a result of this deviation. Clearly then, this is not a best response. Next suppose that an agent deviated by choosing a list  $\ell_i$  such that  $\min \ell_i = \bar{x} > \hat{x}$ . It is immediate that agent i will find a successful match the first round, and will get a payoff of  $F(\bar{x}, \hat{x}) < F(\hat{x}, \hat{x})$ , since he is more cooperative than all the other agents. But he could have gotten payoff  $F(\hat{x}, \hat{x})$  by not deviating, so being more cooperative just gives up payoff and so is not a best response. Finally, suppose that an agent deviated by choosing a list  $\ell_i$  such that  $\min \ell_i = \hat{x}$ . Clearly, he would play  $\hat{x}$  in any stage game and so will match in the first round and get a payoff of  $F(\hat{x}, \hat{x})$ . Since this is the same as he gets by not deviating, this deviation does not improve his welfare.

We conclude that no deviation is a better response, and so the strategy profile given in the hypothesis is an SPE.  $\Box$ 

What this says is that when search is relatively cheap (so agents can easily reject undesirable partners) all agents learning a list consisting of the same unique action is an SPE. We think of this as a social norm since all agents choose to do exactly the same thing and are capable of doing nothing else in equilibrium.

Next we show that other, more complicated lists which lead to the same outcome given above may also be seen in an SPE. Specifically, in equilibrium agents may learn many stage game actions, as long as they have the same lowest (least cooperative) element. This lowest element will then be the only action played in the stage game and so will constitute the social norm.

**Corollary 1.1.** If z is sufficiently small, then for any  $\hat{x} \in \mathcal{X}$  there exists an SPE in which  $\min \ell_i = \hat{x}$  for all  $i \in \mathcal{I}$ , and  $\hat{x}$  is played by all agents in equilibrium in the stage game.

**Proof.** Suppose two agents with these lists were matched. If they agree to play, then in an SPE they will always play the least cooperative strategy on their list since this is the only Nash equilibrium in the subgame. Given that this is the best they could do with any rematching, it is a best response to agree to play rather than to wait to play with another agent in a future round and pay a delay cost. Given this, the addition of other more cooperative strategies to an agent's list is irrelevant in equilibrium. The remaining details of why no deviation from this improves the welfare of any agent follows the same intuition given in the proof of Theorem 1.

The next result states that the noncooperative outcome is always an equilibrium regardless of the delay cost z.

**Theorem 2.** There exists an SPE in which  $\ell_i = \{0\}$  for all  $i \in \mathcal{I}$ , and 0 is played by all agents in equilibrium in the stage game.

**Proof.** Consider a strategy profile in which all agents learn only strategy 0. As above, it is a best response to play this in any stage game and since this is true for all agents, it is a best response for the agent to play with any agent with whom he is matched. Thus, the best any agent can do is take the payoff F(0,0) in the first round. Omitting strategy 0 from your list just gives away payoff (since now you cooperate more than other agents), and adding other more cooperative strategies to your list is irrelevant since they will never be played in equilibrium when 0 is available. It follows that this is an SPE.  $\Box$ 

That social norms are equilibria is perhaps not hugely surprising. The more interesting and difficult part of this work is to show that these are the *only* equilibria. This contrasts with the Folk Theorem for repeated games.

We begin by focusing on a few particular types of strategies and then show that these are the only strategies that will be played in equilibrium. Specifically, we will categorize a player's strategies as having two identifying characteristics:

the least cooperative strategy in their own list  $(x^m = \min\{\ell_i\})$  where the superscript m is mnemonic for "my strategy") and the least cooperative strategy that they will agree to play against  $(x^y)$  where the superscript y is a mnemonic for "your strategy"). More formally, define the set of players using such strategies for any given profile S as follows:

$$\mathcal{I}_{x^m,x^y}(S) \equiv \left\{ i \in \mathcal{I} \mid x^m = \min\{\ell_i\} \text{ and } \forall \ell \in \mathcal{L}, \text{ if } \min\{\ell\} \geqslant x^y, \ g_i(\ell) = x^m, \right.$$
and if  $\min\{\ell\} < x^y, \ g_i(\ell) = N \right\}.$ 

In reading this definition, note that we embed the assumption that if a player i is willing to play against a list  $\ell$  with least element  $x^y$ , then i is also willing to play against any list with a larger least element. This is immediate since if it is a best response to play with a less cooperative player it is also a best response to offer to play with a more cooperative player. In addition, we assume that players always employ the least cooperative strategy in their own lists when they agree to play with any opponent. Again, it is immediate that this is the only subgame perfect best response. Thus, if a strategy is part of an SPE, it must fall into some category

$$\mathcal{I}_{x^m,x^y}(S)$$
.

We can therefore restrict our attention to such strategies from this point forward.

Now that we know the form that equilibrium strategies must take, it is useful to form three exhaustive subcategories. It is possible that agents insist that their opponents behave in a more cooperative way than they are willing to themselves. We will call these exploitive agents *wolves*. On the other hand, it may be that agents are willing to accept a match with an opponent who is less cooperative than they are, perhaps in order to avoid the cost of continued searching. We call these more exploitable agents *sheep*. Finally, agents may be willing to play only with agents who are at least as cooperative as they are willing to be. We call agents who insist on being treated at least symmetrically *norm players*. We will add a superscript to the partition we defined above to keep track of this:

```
\mathcal{I}_{x^m, x^y}(S) \subset \mathcal{I}^w(S) if and only if x^m < x^y, \mathcal{I}_{x^m, x^y}(S) \subset \mathcal{I}^s(S) if and only if x^m > x^y, \mathcal{I}_{x^m, x^y}(S) \subset \mathcal{I}^n(S) if and only if x^m = x^y.
```

We will call wolves and sheep *complementary* players if the wolves choose to behave in the least cooperative way the sheep will accept, and the sheep behave in the least cooperative way the wolves will accept. More precisely, consider  $\mathcal{I}_{x^{mw},x^{yw}}(S) \subset \mathcal{I}^w(S)$  and  $\mathcal{I}_{x^{ms},x^{ys}}(S) \subset \mathcal{I}^s(S)$ . Then these agents are said to have *complementary strategies* if and only if  $x^{mw} = x^{ys}$  and  $x^{ms} = x^{yw}$ .

Note that the equilibria shown to exist in Theorems 1 and 2 involve all agents in the game playing the same strategies. Although it might be that agents use slightly different strategies to support this outcome (for example, learning different lists, but with the same least element) at any SPE, the outcome is that all agents abide by the social norm.

It turns out there is one other type of equilibrium which also has the flavor of a social norm. It is possible for all the agents but one to choose to be wolves, and for the one remaining agent to be the complementary type of sheep when search costs are small. Being a wolf is almost a social norm in this case. This equilibrium is not robust to additional refinement and to this extent is an artifact of the equilibrium concept (in contrast to the more realistic social norm SPE discussed above). While the one sheep in this equilibrium cannot improve his payoff by defecting from this strategy in a static sense, any reasonable sheep would see that if he chooses to become a norm player, for example, the wolves would have no choice but to follow his example. We do not consider this kind of far-sighted behavior in this paper, however, and so the equilibrium remains.

**Theorem 3.** For small enough z, there exists an SPE in which there is exactly one sheep, and all other agents are complementary wolves.

**Proof.** With so many wolves hoping to find the one sheep, wolves must search many periods before finding a match. Sheep find a match right away, of course. As z becomes small, however, the net search costs paid by wolves goes to zero, and the expected payoff converges to the stage game payoff. Given this, no wolf would be better off becoming a sheep since this results in a lower payoff. Also, adopting any other strategy that results in the sheep agreeing to play when matched only gives up payoff since the complementary wolf strategy results in the largest possible payoff in this

case. The only alternative strategies that would result in a wolf agent being able to play when matched are to choose to instead to become a sheep even more submissive than the one that already exists. Clearly, this will not increase his payoff. Thus, it is a best response for a wolf to remain a wolf.

Why doesn't the sheep choose a different strategy? Clearly he won't choose to be a more submissive sheep as this only lowers his payoff. Choosing to be a more aggressive sheep, a norm player or a wolf, however means that neither he nor any other player ever finds a successful match since all the other players match only with complementary sheep or sheep that are even more submissive. Thus, defecting from the original sheep strategy results in a payoff of -z in each period. Clearly this is worse than accepting whatever payoff he gets from being the one lone sheep.  $\Box$ 

At last we are ready to show that the only SPE are those discussed in the propositions above. This is done through a series of lemmas given in Appendix A that lead up to the following theorem.

**Theorem 4.** For z small enough, only two types of SPE are possible: (1) all agents using the same norm strategy, (2) one sheep player and all other agents using the complementary wolf strategy.

## **Proof.** See Appendix A. $\Box$

Note that if we set  $\epsilon = 1$ , we get back the standard pure strategy prisoners' dilemma game as a special case. The theorems then say that all agents choosing to learn c only and all agents choosing to learn d only (or equivalently, c and d then playing d) are the only two robust SPE.

Another immediate corollary is the following:

**Corollary 4.1.** For any given z if  $\epsilon$ , the gap between pure strategies, is small enough, then the only equilibrium outcome is for all agents to learn to defect (x = 0) and to play this in every encounter.

**Proof.** By Theorem 2, the noncooperative outcome is always an SPE. Thus, we need only show this is the only SPE for small  $\epsilon$ . Suppose that an agent i had a strategy g that required that his opponent j employ strategy x > 0 and threatened to pass otherwise. Suppose that the opponent j's least cooperative strategy was  $x - \epsilon$ . Since for small  $\epsilon$ , cost of delay z exceeds the small  $(\epsilon)$  loss from playing with a slightly less cooperative opponent this period, it is not credible (under the definition of SPE) for agent i to decline to play with such an agent j. It follows that no such strategy (requiring that an opponent j employ any specific strategy x > 0 and threatening to pass otherwise) could be subgame perfect. The only remaining strategy is to learn to never cooperate (x = 0) and play this against all partners in all rounds.  $\square$ 

Intuitively, this says the following: suppose that we begin at any partially cooperative social norm SPE. Clearly, any one agent would improve his payoff by being just a little bit less cooperative. Provided he does not take away more from his potential partner than the search costs, he will still find a match in the first round. But then all agents are better served by learning a slightly less cooperative strategy and this is the social norm. Of course, then it is optimal for any given agent to be slightly less cooperative than this new social norm. This process unravels until only the fully noncooperative social norm remains. An immediate implication of this is that if the strategy set is continuous (or if mixed strategies are allowed), the game unravels and only the noncooperative strategy remains as an equilibrium

This suggests that for any kind of cooperation to be possible, there must be clear lines between ethical systems. To be seen in equilibrium, a philosophy or creed must lay out a clear code of behavior and must be measurably distinct from alternative ethical codes.

The assumption of evenly spaced granularity between strategies could also be relaxed. Suppose that there is an asymmetry in the gap between ethical systems. For example, once one falls into an ethical muddle of cheating with any frequency, it might be hard to think of any ethical system that would permit cheating only to one particular degree. On the other hand, one can imagine a system that says never cheat. If this is the case, then once you fall below a certain ethical standard of behavior, there is not sufficient granularity to prevent you from falling to the bottom of the ethical scale. More formally, given search costs z, if the gap between strategies becomes smaller as one moves toward pure noncooperation then the set of potential SPEs will include all symmetric outcomes above some cut point at which the granularity become too fine (and as always, the noncooperative outcome as well).

#### 4. The anti-prisoner's dilemma and harmful social norms

| Anti-prisoner's dilemma |        |           |  |
|-------------------------|--------|-----------|--|
|                         | Behave | Misbehave |  |
| Behave                  | 10, 10 | 12, 0     |  |
| Misbehave               | 0, 12  | 2, 2      |  |

While much has been written (including this paper) about the prisoner's dilemma, little attention has been devoted to games like the anti-prisoner's dilemma depicted above. It is easy to see why. The game has a dominant strategy equilibrium which is Pareto efficient. Equilibrium play in both the one-shot and repeated game case will result in a Pareto efficient outcome. Thus, the game appears trivial, and in particular to present the players no dilemma at all. We will see below, however, that this game yields very different results when placed in an endogenous game structure.

To give this example more context, suppose that high school kids form groups of two (best friends) to do things together. They can choose to be good kids, work hard at school, and generally behave, or skip classes, get in trouble and generally be delinquent. The payoff from behaving is being allowed privileges such as a car, late curfew, and praise from one's parents. Misbehaving results in these privileges being taken away, encounters with the police and disapproval from one's parents. In our example, the rewards from behaving exceed whatever pleasures there are from misbehaving. Parents, however, tend to judge their children by comparing them to their peers. Thus, if a child is misbehaving much more than his best friend, punishments and disapproval are more severe. Conversely, if a child behaves well compared to his friends, he looks especially good and gets exceptional rewards. From the prospective of the kids, behaving is a dominant strategy. The best possible outcome, however, is to be a very well behaved kid with an extremely delinquent friend. This maximizes the good kid's payoff (and in turn, minimizes that of the bad kid).

Considered as a one-shot game both kids behaving is the dominant strategy equilibrium. However, if the game is played endogenously (and the waiting cost is sufficiently small) there exists a subgame perfect equilibrium in which everyone misbehaves. To see this, suppose that all kids know only how to "misbehave." If such kids choose to play with a kid who knows how to behave, they end up looking rotten and getting the worst payoff possible. "Bad kids" are better off waiting for another bad kid to come along. As a result, no kid would ever want to find a friend who has learned how to behave in the pregame and so kids know that if they choose to be "good" they will be ostracized. Thus, the endogenous game leads to a social norm in which all kids know only how to misbehave and so everyone obtains a payoff below the minmax payoff of 10.9

The endogenous anti-prisoner's dilemma captures a particularly harsh form of peer pressure. No one likes misbehaving in our example. Nevertheless, everyone does so in order to "fit in." Thus, we can see how social norms enforced through ostracism can be harmful to all.

We conclude this section with one final point. It may seem from the above examples that any symmetric outcome of a symmetric game can be supported as the subgame perfect equilibrium of an endogenous game. This is not the case, as can be seen in the following game:

| Game with a unique equilibrium |        |      |  |
|--------------------------------|--------|------|--|
|                                | A      | В    |  |
| A                              | 10, 10 | 3, 4 |  |
| В                              | 4, 3   | 2, 2 |  |

In this game there is no equilibrium in which agents have only B in their lists, no matter how small waiting costs are. To see why, suppose that everyone has only B in their lists, and that agent i deviates during the pregame and chooses the list A. When agent i is matched with agent j, agent j must opt to play against agent i because agent j earns 3 from playing against i but earns 2 from playing against anyone else. Consequently, this endogenous game has a unique class of subgame perfect equilibria in which every list contains action the dominant action A, and the only equilibrium payoff is the minmax payoff of 10.

<sup>&</sup>lt;sup>9</sup> Of course if we allowed kids differing degrees of possible misbehavior (that is, granular strategies) all symmetric social norms of partial misbehavior would also be equilibria. It is also the case that these are the only SPE, just as in the standard prisoner's dilemma.

#### 5. Literature review

The starting point for this work is Rubinstein and Wolinsky's (1985) famous bargaining model. This began a large literature too extensive to be cited in detail here. <sup>10</sup> Their paper provided a very useful formalization of the notion that the players in a game are not randomly determined, but coalesce endogenously instead.

Although Rubinstein and Wolinsky do not further consider that the form of the game that agents play may also be endogenously determined, there are a number of other papers that do so in various ways. An excellent recent example is Jackson and Wilkie (2005). The games they consider are endogenous in the sense that agents are able to enter into certain classes of binding contracts to give one another side payments contingent on the strategies they eventually choose. This is an example of precommitment being introduced into the game and is in the spirit of Coasian contracting. It has the effect of altering the payoff function for the agents, but not the strategy sets they have available. There are many variations on this kind of pregame contracting, see especially, Varian (1994), Ray and Vohra (1997), Caruana and Einav (2008) and other work discussed in Jackson and Wilkie (2005).

A somewhat different approach is exemplified by Lagunoff (1992). This literature imagines that there is an initial phase in which the overall mechanism that agents will use to decide on allocations is chosen. Thus, the entire form of the game is up for discussion. Related to this is a literature that studies the optimal size of the agenda in bargaining games. Lang and Rosenthal (2001), for example, explore whether it is beneficial to bargain over a large agenda simultaneously or only over a subset. Again, both of these papers and the literatures they represent consider macrolevel alterations of the whole game form through planning, voting, or negotiation. This is different from the micro scale approach proposed here of allowing agents to voluntarily and unilaterally restrict their strategy choices within a given game framework.

Alger and Renault (2006, 2007) examine a principal-agent model in which agents' types are two-dimensional: they have an ethical type and a cost type. Honest agents cannot lie about their cost type but they can lie about their honesty, while dishonest agents can lie about both their ethical type and their cost type. The papers explore the ability of the principal to screen agents according to ethical type when agents first send messages about ethical type and then send messages about their cost type. Alger and Renault's model is similar to ours in that an agent's type is related to whether the principal wants to transact with him and restricts the transactions that can take place. The principal must design a contract to induce an honest agent to reveal his restricted strategy set. In our model the agent must reveal the restricted strategy set or, put differently, the list is also the only signal available to the agent.<sup>12</sup>

Gilson and Mnookin (1994) present an informal model in which litigants can choose lawyers with reputations for cooperation, and thereby commit to a more cooperative litigation procedure than they would have been able to achieve otherwise. If one considers choosing a lawyer as choosing a list, their informal model would be a special case of ours. Croson and Mnookin (1997) find experimental support for the Gilson and Mnookin model.

#### 6. Remarks

**Remark 1** (*Precommitment*). There are several ways to model precommitment. The most typical is for agents (monetary authorities, for example) to choose policies that are irreversible in dynamic games. Another is contingent contracts that assure your opponent that you will play in a certain way due to an *ex post* enforceable penalty clause (most favored customer agreements, for example). In contrast, agents in our game have the option of narrowing the set of strategic options open to themselves and do so in response to the choices made by the population of agents at large. Contingent contracts are not allowed and the potential equilibrium outcomes are, in principle, richer than could be obtained from requiring that agents choose only a single strategy for all time at the beginning of the game. The results are also different: in some cases, many different social norms can emerge as equilibrium outcomes while in others (high search costs) only the one-shot Nash equilibrium will be seen despite the possibility of altering the game. Perhaps most importantly, committing not to play certain strategies in the pregame is only optimal because agents can

 $<sup>^{10}</sup>$  See, for example, Osborne and Rubinstein (1990) as well as the other work discussed in this section.

<sup>&</sup>lt;sup>11</sup> Bade et al. (2006) allow agents to commit to a subset of actions, similar to our lists. In their paper, though, lists only restrict choice while here lists both restrict choice and form the basis on which matches are consummated.

<sup>&</sup>lt;sup>12</sup> In related work, Chen (2000) examines whether incomplete contracts can be socially preferable when a subset of agents have strategy sets that restrict them from breaking promises.

also choose their opponents/partners. No agent would choose to limit his strategic options in the prisoners dilemma if he did not have the option of declining to play with an aggressive opponent. Thus, while the approach we take does have an element of precommitment, the form it takes, the elements that inform its use, and the results it yields are different from what is typically seen in the literature.

Remark 2 (Learning, ethics and observability). We motivate the pregame in two ways which we see as independently interesting. One reason that strategies might be unavailable to me is that I may have never learned them. For example, telling lies to cover one's tracks takes practice and learning how to attract the opposite sex is difficult. Whether or not an agent has learned certain skills is, at least to an extent, observable. For example I might infer from your accent or word choice you can speak Spanish. Thus, in some contexts learning makes the commitment not to employ a strategy credible and also contributes to the observability of these commitments. We might instead motivate removing options from a player's strategy set as reflecting his ethical framework. This is a kind of axiomatic approach to ethics with clear lines drawn between what is permitted and what is prohibited. Again, there are frequently outward signs of one's choice of ethical framework and reasons to believe that agents will not deviate from these internal standards of behavior. For example, observing an agent in dress of an orthodox Jew allows us to infer something about his code of behavior. Of course, one could alternatively treat excluded strategies as having large negative payoffs instead. It might very well be that the outcomes would be the same under this alternative formation, but it seems more natural to us to treat things you have never learned or you would never consider doing as simply excluded from your strategy set. Finally, there are certainly examples of real world prisoners' dilemmas and other games in which it would be very hard to argue that ethics could constrain you from paying a strategy or that any meaningful effort has to be put in to learn a strategy or that any of this could be observable in any event. An example might be checkers. In this case, our model would simply not apply. While we think the endogenous game we develop in this paper has wide application, we certainly are making no claim of universal applicability.

#### 7. Conclusion

We describe a research program that examines a new form of game in which both players and the exact form of the game (strategy sets) are endogenously determined through equilibrium play. We argue that such endogenous games are not uncommon in the real world, and indeed may be more the rule than the exception. We explore the granular prisoners' dilemma and anti-prisoners' dilemma in detail. We show that if the search cost is low enough, the only robust equilibria can be seen as social norm outcomes in which all agents cooperate to exactly the same degree.

We make a number of simplifying assumptions in showing these results and it would be interesting to relax these. For example, we assumed that agents can perfectly observe the list of the agent with whom he is matched. One might instead include costly signaling about the lists agents choose, or perhaps probabilistic revelation of lists as in Dekel et al. (2007). We also assumed that it was costless to learn strategies. While making learning costly in the PD games we explore here would not change anything, in more general games it might make a difference. Finally, we could allow agents who replace those who match out to choose their own strategies. We think this will lead to somewhat different results, but not ones that give much new insight. For example, in the granular PD games, all agents match the period they are born and so we could see each new generation settling on a new and different social norm. This will dramatically increase the complexity of the equilibrium notion and strictly increase the number of equilibrium outcomes. Nevertheless, we still see social norms arising in each period.

This work could be extended in several directions. Most obviously, one could consider other types of games. In separate work, we explore coordination and ultimatum games. One could also move away from specific games to see what could be proven about endogenous games in general. Finally, these games have a very simple structure and it would be quite interesting to see if our theoretical results hold up to experimental verification.

For the present, our point is to argue that if we have to convince people to play games with us we might find it in our interests to change the sort of person we are. We think this sheds new light on the important findings of the experimental and behavioral literatures. It may be that agents are in fact behaving in a fully rational way given the strategy sets they have available to them. The findings of these literatures may therefore be seen as perhaps rejecting the embedded assumption that agents are playing a specific game rather than that they are playing not in a fully rational way.

## Appendix A

In this appendix, we prove a series of lemmas that lead up to showing that essentially nothing but norm strategies can be an SPE. We begin by showing that only one type of norm player can appear in equilibrium

**Lemma 1.** It is impossible for two types of social norm players to coexist in any SPE.

**Proof.** Suppose there were two nonempty sets of norm players  $\mathcal{I}_{\hat{x},\hat{x}}(S)$  and  $\mathcal{I}_{\bar{x},\bar{x}}(S)$ ,  $\hat{x} \neq \bar{x}$ , at an SPE. Without loss of generality, suppose the expected payoff players in the set  $\mathcal{I}_{\hat{x},\hat{x}}(S)$  got was at least as big as the expected payoff to players in  $\mathcal{I}_{\bar{x},\bar{x}}(S)$ . Suppose one agent  $i \in \mathcal{I}_{\bar{x},\bar{x}}(S)$  chooses to become the other kind of norm player. The payoff to agents in  $\mathcal{I}_{\hat{x},\hat{x}}(S)$  would go up since the addition of one new member would lower the expected number of periods needed to find a partner with whom to play. Thus, it is not a best response for agent i to choose to be a  $\bar{x}$  norm player instead of a  $\hat{x}$  norm player and so this could not have been an SPE.

Next we show that if any wolf players exist in an SPE, there must also exist at least one sheep player.

Lemma 2. If there exists any wolf player in an SPE, then there must also exist at least one sheep player.

**Proof.** Observe that wolf players will never play the stage game when paired with any norm or other wolf player. If paired with a norm player who would agree to play with him, it must be that the wolf player will choose to pass, and when paired with a norm player with whom the wolf player would play the norm player would always pass. Similarly, if two wolf players are paired, one of them will find that the other is not cooperative enough and so will pass. It follows that if there are no sheep, the wolf players will fail to play each period and so get a payoff of -z each period. On the other hand, by defecting and becoming a complementary sheep, the worst they can do is get a one-time payoff of -3. This is clearly a higher payoff so choosing to be a wolf is not a best response when there are no sheep.  $\Box$ 

Lemma 3 shows that norm and sheep players will not both appear in equilibrium when search costs are small enough.

**Lemma 3.** For z small enough, it is impossible for sheep and norm players to coexist in any SPE.

**Proof.** By Lemma 1, there can only exist one type of norm player at an SPE. We consider two subcases.

First, suppose that the norm player and some kind of sheep player have strategies that result in their agreeing to play when matched. Then norm players get a weakly higher payoff if they are lucky enough to be matched with this kind of sheep compared to what they receive when matched to another norm player. This is because such sheep play stage game strategies that are at least as cooperative as the ones used by the norm players. If the payoff is strictly higher, then for z small enough, waiting for these sheep becomes less costly and it ultimately becomes more profitable for a norm player to wait for these sheep than to agree to play with another norm player. Therefore, it is a best response for any norm player to defect and become a complementary wolf. Observe that this is also true if the payoff happens to be the same. In this case, sheep are playing exactly the same stage game strategy as the norm players but would also play with a less cooperative wolf. Thus, by defecting and becoming this complementary wolf, the payoff the defecting norm player receives is also higher when matched with a sheep than playing the norm strategy with a sheep. As a result, the same logic applies

Second, suppose instead that for some kind of sheep appearing in the SPE, the sheep and norm players won't agree to play together. Then for any z, either the norm player payoff is at least as high as the sheep player payoff or inversely. If the norm payoff is larger, then by becoming a norm player, the sheep player increases the expected payoff of norm players, himself included, since this reduces the time to find a match while not changing the expected payoff to a successful match. If the sheep players get at least as high a payoff as norm players, then opposite is true. By becoming a sheep the expected time to find a match for all sheep goes down, while the expected payoff to successful matches goes up.

Thus, sheep and norm players cannot coexist for small enough z.  $\Box$ 

Next we show that wolves and norm player will not be found together in an SPE for small enough search costs.

**Lemma 4.** For z small enough, it is impossible for wolves and norm players to coexist in any SPE.

**Proof.** By Lemma 2, if wolves exist at an SPE then sheep also exist for small enough z. But if sheep exist then by Lemma 3, no norm players can exist in an SPE for small enough z. The conclusion follows.  $\Box$ 

The next lemma shows that as search costs get small, in the limit at most one sheep can remain in any SPE. All the others will choose to defect and choose other strategies.

**Lemma 5.** For z small enough, at most there will be one sheep agent in any SPE.

**Proof.** Suppose there existed some type of sheep with at least two representatives in an SPE for all z. All the other agents in the game must be either wolves or other sheep by Lemma 3.

Suppose that z gets very small. Then the payoff from being the complementary kind of wolf to this sheep approaches the stage game payoff (since the expected search costs become zero in the limit). On the other hand, the expected payoff to being this kind of sheep are strictly lower than the wolf payoffs since in all cases the sheep play a more cooperative strategy when matched. It follows that if z is small enough, regardless of what strategies the other players in the game choose, if there are least two of these sheep, it is a best response for one of them to defect and become the complementary wolf. (Note this does not hold if there is only one agent who chooses a given sheep strategy since by defecting, he obliterates his entire food supply. Also note that for small enough z, the wolves will get a strictly higher payoff than the sheep, but the one remaining sheep still cannot profitably defect because of the obliteration issue above.)

Now suppose that there are two agents who choose to be different kinds of sheep. The same logic holds. At least one of these agents will be better off defecting and becoming a complementary wolf to the other sheep for small enough z.

We conclude there will be at most one sheep at any SPE for small enough z.  $\Box$ 

Our last lemma shows that when search costs are small, we will see at most one kind of wolf in any SPE.

**Lemma 6.** For z small enough, at most there will be one kind of wolf at any SPE.

**Proof.** By Lemmas 4 and 5, if wolves exist at an SPE, there can be no norm players and at most one sheep for small z. If there is only one sheep, all wolves must be of the complementary type. If they are too greedy and cannot match with the lone sheep, they are better off becoming complementary wolves since then they get a positive payoff instead of -z forever. If they are too generous and give away more than is necessary play with the sheep, they should defect and become complementary wolves and get a higher payoff. Thus, only complementary wolves will exist in equilibrium.  $\Box$ 

Putting these lemmas together allows us to prove that the types of equilibria shown to exist in the body of the paper are in fact the only SPE that exist.

**Theorem 4.** For z small enough, only two types of SPE are possible: (1) all agents using the same norm strategy, (2) one sheep player and all other agents using the complementary wolf strategies.

**Proof.** Since we know that wolf, sheep and norm strategies exhaustively partition the set of strategies that will be used in any SPE, by Lemma 1, at most one type of norm strategy can appear in any SPE. By Lemmas 3 and 4, norm players will not be seen with either sheep or wolf players in equilibrium when search costs are small enough. Lemmas 5 and 6 tell us that for small enough z, there can be at most one kind of wolf and one individual sheep player in any SPE. This means that the only configurations possible for an SPE if z is small enough are (1) all agents using the same norm strategy, and (2) one sheep player and all other agents using identical wolf strategies. Note that for (2), the wolf strategy must be complementary to the one employed by the sheep since this is only best response.

#### References

Alger, Ingela, Renault, Regis, 2006. Screening ethics when honest agents care about fairness. Int. Econ. Rev. 47, 59-85.

Alger, Ingela, Renault, Regis, 2007. Screening ethics when honest agents keep their word. Econ. Theory 30, 291-311.

Axelrod, Robert, 1986. An evolutionary approach to norms. Amer. Polit. Sci. Rev. 80, 1095-1111.

Bade, Sophie, Haeringer, Guillaume, Renou, Ludovic, 2006. Bilateral commitment. Working Paper. Penn State University.

Baker, Matthew J., Jacobsen, Joyce P., 2007. A human capital-based theory of postmarital residence rules. J. Law, Econ., Organ. 23, 209-241.

Bendor, Jonathan, Swistak, Piotr, 2001. The evolution of norms. Amer. J. Sociology 106, 1493-1545.

Burdett, Ken, Coles, Melvyn G., 2001. Transplants and implants: The economics of self-improvement. Int. Econ. Rev. 42, 597-616.

Caruana, Guillermo, Einav, Liran, 2008. A theory of endogenous commitment. Rev. Econ. Stud. 75, 99-116.

Chen, Yongmin, 2000. Promises, trust, and contracts. J. Law, Econ., Organ. 16, 209–232.

Coleman, James S., 1990. Foundations of Social Theory. Harvard Univ. Press.

Cole, Harold L., Mailath, George J., Postlewaite, Andrew, 2001. Efficient non-contractible investments in large economics. J. Econ. Theory 101, 333–373.

Croson, Rachel, Mnookin, Robert H., 1997. Does disputing through agents enhance cooperation? Experimental evidence. J. Legal Stud. 26, 331–345.

Dal Bo, Pedro, 2007. Social norms, cooperation and inequality. Econ. Theory 30, 89-105.

Dekel, Eddie, Ely, Jeffrey C., Yilankaya, Okan, 2007. Evolution of preferences. Rev. Econ. Stud. 74, 685-704.

Ellison, Glenn, 1994. Cooperation in the prisoners dilemma with anonymous random matching. Rev. Econ. Stud. 61, 567–588.

Gilson, Ronald J., Mnookin, Robert H., 1994. Disputing through agents: Cooperation and conflict between lawyers in litigation. Columbia Law Rev. 94, 509–549.

Hirshleifer, David, Rasmusen, Eric, 1989. Cooperation in a repeated prisoners dilemma with ostracism. J. Econ. Behav. Organ. 12, 87-106.

Jackson, Matthew, Wilkie, Simon, 2005. Endogenous games and mechanisms: Side payments among players. Rev. Econ. Stud. 72, 543-566.

Jansen, Marcel, 2004. Can job competition prevent hold-ups?. IZA Discussion Paper No. 988.

Kandori, Michihiro, 1992. Social norms and community enforcement. Rev. Econ. Stud. 59, 63-80.

Lagunoff, Roger, 1992. Fully endogenous mechanism selection on finite outcome sets. Econ. Theory 2, 462–480.

Lang, Kevin, Rosenthal, Robert, 2001. Bargaining piecemeal or all at once. Econ. J. 111, 526-540.

Neilson, William S., 1999. The economics of favors. J. Econ. Behav. Organ. 39, 387-397.

Osborne, Martin J., Rubinstein, Ariel, 1990. Bargaining and Markets. Academic Press.

Peters, Michael, Siow, Aloysius, 2002. Competing premarital investments. J. Polit. Economy 110, 592-608.

Ray, Debraj, Vohra, Rajiv, 1997. Equilibrium binding agreements. J. Econ. Theory 73, 30-78.

Rogerson, Richard, Shimer, Robert, Wright, Randall, 2005. Search-theoretic models of the labor market: A survey. J. Econ. Lit. 43, 959-988.

Rubinstein, Ariel, Wolinsky, Asher, 1985. Equilibrium in a market with sequential bargaining. Econometrica 53, 1133-1150.

Sugden, Robert, 2004. The Economics of Rights, Co-Operation and Welfare. Palgrave Macmillan.

Varian, H., 1994. A solution to the problem of externalities when agents are well-informed. Amer. Econ. Rev. 84, 1278–1293.