

ok; add 10m t

# Reputation and Equilibrium Characterization in Repeated Games with Conflicting Interests

by Klaus M. Schmidt<sup>1</sup>

Department of Economics, E52-252C  
Massachusetts Institute of Technology  
Cambridge, MA 02139, USA  
Tel: (617) 253-6217

March 1991

revised: February, 1992

A two-person game is of *conflicting interests* if the strategy to which player one would most like to commit herself holds player two down to his minmax payoff. Suppose there is a positive prior probability that player one is a “commitment type” who will always play this strategy. Then player one will get at least her commitment payoff in any Nash equilibrium of the repeated game if her discount factor approaches one. This result is robust against further perturbations of the informational structure and in striking contrast to the message of the Folk theorem for games with incomplete information.

KEYWORDS: Commitment, Folk theorem, Repeated Games, Reputation.

HEADNOTE: Repeated Games with Conflicting Interests

---

<sup>1</sup>This paper is based on Chapter 3 of my PhD thesis which was completed within the European Doctoral Programme at Bonn University. I would like to thank David Canning, In-Koo Cho, Benny Moldovanu, Georg Nöldeke, Ariel Rubinstein, Avner Shaked, Joel Sobel, Monika Schnitzer, Eric van Damme and in particular Drew Fudenberg for many helpful comments and discussions. Financial support by Deutsche Forschungsgemeinschaft, SFB 303 at Bonn University, is gratefully acknowledged.

# 1 Introduction

Consider a repeated relationship between two long-run players one of whom has some private information about her type. A common intuition is that the informed player may take advantage of the uncertainty of her opponent and enforce an outcome more favourable to her than what she would have got under complete information. This intuition has been called “reputation effect” and has found considerable attention in the literature. The purpose of this paper is to formalize this intuition in a general model of repeated games with “conflicting interests” and to show that the effect is robust against perturbations of the informational structure of the game.

The first formalizations of the “reputation effect” in games with incomplete information have been developed by Kreps and Wilson (1982) and Milgrom and Roberts (1982). They have shown that a small amount of incomplete information can be sufficient to overcome Selten’s (1978) chain-store paradox. An incumbent monopolist who faces a sequence of potential entrants may deter entry by maintaining a reputation for “toughness” if there is a small prior probability that she is a “tough” type who prefers a price war to acquiescence. Recently, this result has been generalized and considerably strengthened by Fudenberg and Levine (1989, 1992). They consider the class of all repeated games in which a long-run player faces a sequence of short-run opponents, each of whom plays only once but observes all previous play. They show that if there is a positive prior probability of a “commitment type”, who always plays the strategy to which player one would most like to commit herself, and if player one is sufficiently patient, then she can enforce at least her commitment payoff in any Nash equilibrium, i.e. she will get at least what she would have obtained if she could have committed herself publicly to this strategy. This result is very powerful, because (i) it gives a tight lower bound for player one’s payoff in all Nash equilibria, (ii) it holds for finitely and infinitely repeated games, and (iii) it is robust against further perturbations of the informational structure, i.e. it is independent of what other types may exist with positive probability. However, Fudenberg and Levine’s analysis is restricted to games where a long-run player faces a sequence of short-run opponents. Our paper provides a generalization and qualification of their results for the two long-run player case. We show that a necessary and sufficient condition for this generalization to hold is that the game is of conflicting interests.

The first who studied reputation effects in repeated games with two long-run players were again Kreps, Wilson, Milgrom and Roberts (1982). They demonstrated that in all sequential equilibria of the finitely repeated prisoner's dilemma the cooperative outcome will be achieved if there is a small prior probability that the players are of a type which is committed to always play the "tit-for-tat" strategy. In the beginning of the 1980s repeated games with incomplete information have been very popular to solve all kind of puzzles in game theory, industrial organization, macroeconomics etc.. But the initial enthusiasm has been considerably dampened by a Folk-theorem type result of Fudenberg and Maskin (1986). They have shown that for any finitely or infinitely repeated game there exists an  $\epsilon$ -perturbation of this game (in which each of the players has a different payoff function with an arbitrarily small but positive prior probability) such that any individually rational, feasible payoff vector of the unperturbed game can be sustained as the outcome of a sequential equilibrium of the perturbed game, if the players are sufficiently patient and if there are enough repetitions. But if any outcome can be "explained" by just picking the right perturbation then the predictive power of the theory of repeated games is very limited indeed.

However, before this message is accepted, we should have a closer look at what Fudenberg and Maskin mean by an  $\epsilon$ -perturbation. They assume that with an arbitrarily small but positive probability each of the players may be "crazy", i.e. she may have a completely different payoff function as compared to the original game. For any given payoff vector Fudenberg and Maskin then pick one very specific type of "crazyness" which is used to sustain this payoff vector as an equilibrium outcome. A perturbation of the informational structure of a game should capture the idea that the players may be slightly uncertain about what the exact payoff functions of their opponents are. This may include the possibility that the opponent is actually crazy as in Fudenberg and Maskin (1986), but it should not be restricted to only one (very particular) type of crazyness. We argue that the modeller should allow for a broader class of perturbations, such that many different payoff functions may have positive prior probability. Surprisingly, insisting on this kind of robustness yields a result which is in striking contrast to the message of the Folk-theorem. No matter what types may possibly be drawn by nature (including those considered by Fudenberg and Maskin) and how likely they are to occur, if player one is sufficiently patient, if the game is of "conflicting interests", and if there is an arbitrarily small but

positive probability of a “commitment” type, then we can give a tight prediction of the equilibrium outcome of all Nash equilibria.

To make this more precise, consider a repeated game with complete information in which player one would like to commit herself to take an action  $a_1^*$ , called her “commitment action”, in every period. If player two responds optimally to  $a_1^*$  player one gets her “commitment payoff”. Assume that the game is of “conflicting interests” in the sense that playing  $a_1^*$  holds player two down to his minimax payoff. Now suppose that the informational structure of this game is perturbed such that player one may be one of several possible “types”. Consider a type for whom it is a dominant strategy in the repeated game always to play  $a_1^*$  and call her the “commitment type”. Our main theorem says that if the commitment type has any arbitrarily small but positive probability and if player one’s discount factor goes to one then her payoff in any Nash equilibrium is bounded below by her commitment payoff. This result is independent of the nature of the other possible types and their respective probabilities. We generalize the theorem to the case of two-sided uncertainty. Furthermore, we show that “conflicting interests” are a necessary condition for our theorem to hold.

Our result highlights the importance of the relative patience of the two players. Player one has to be sufficiently patient as compared to player two, i.e. for any given discount factor  $\delta_2 < 1$  there exists a  $\underline{\delta}_1(\delta_2) < 1$  such that player one can enforce his commitment payoff in all Nash equilibria if his discount factor satisfies  $\delta_1 > \underline{\delta}_1(\delta_2)$ . The importance of the relative patience of the two players is most intuitive in the case of a completely symmetric game with two-sided uncertainty. If this game has conflicting interests, then it is clearly not possible that both players get their most preferred outcomes at the same time. However, if one of them is sufficiently more patient (or if the prior probability that she is the commitment type is sufficiently higher) then the reputation effect works to her advantage.

A complementary analysis to ours is Aumann and Sorin (1989). For a different class of repeated games, coordination games with “common interests”, they obtain a similar result. However, they have to restrict the possible perturbations to types who act like automata with bounded recall. They show that if all strategies of recall zero exist with positive probability then all pure strategy equilibria will be close to the cooperative

outcome. In contrast to Aumann and Sorin we allow for any perturbation of player one's payoff function and for mixed strategy equilibria. Games of "common" and of "conflicting" interests are two polar cases. We will discuss them in more detail in Section 5.

Finally, in a very recent paper Cripps and Thomas (1991) characterize the set of Nash equilibria of infinitely repeated games with one-sided incomplete information in which players maximize the limit of the mean of their undiscounted payoffs. Following a different method pioneered by Hart (1985) they also find that in games of conflicting interests the informed player can enforce her commitment payoff if there is an arbitrarily small prior probability of a commitment type. Since there is no discounting their result seems to indicate that the relative patience is not that important after all. However, as we will show in Section 4, this interpretation is misleading.

The rest of the paper is organized as follows. In the next section we introduce the model following closely Fudenberg-Levine (1989) and we briefly summarize their main results. Then we give a counterexample showing that their theorem cannot carry over to the class of all repeated games with two long-run players. This gives some intuition on how this class has to be restricted. Section 4 contains our main results. There we generalize Fudenberg-Levine's (1989) theorem to the two long-run player case, and we show that the restriction to games with "conflicting interests" is a necessary condition for this generalization to hold. Furthermore we extend the analysis to the case of two-sided incomplete information. In Section 5 we give several examples which demonstrate how restrictive the "conflicting interests" condition is. Section 6 concludes and briefly outlines several extensions of the model.

## 2 Description of the Game

In most of the paper we consider the following very simple model of a repeated game which is an adaptation of Fudenberg-Levine (1989) and Fudenberg-Kreps-Maskin (1990) to the two long-run player case. The two players are called "one" (she) and "two" (he). In every period they move simultaneously and choose an action  $a_i$  out of their respective action sets  $A_i$ ,  $i \in \{1, 2\}$ . Here we will assume that the  $A_i$  are finite sets.<sup>2</sup> As a point of

---

<sup>2</sup>See Section 6 for the extension to extensive form stage games, continuous strategy spaces and more than two players.

reference consider the unperturbed game (with complete information) first. Let  $g_i(a_1, a_2)$  denote the payoff function of player  $i$  in the unperturbed stage game  $g$  depending on the actions taken by both players. Let  $\mathcal{A}_i$  denote the set of all mixed strategies  $\alpha_i$  of player  $i$  and (in an abuse of notation)  $g_i(\alpha_1, \alpha_2)$  the expected stage game payoffs.

The  $T$ -fold repetition of the stage game  $g$  is denoted by  $G^T$ , where  $T$  may be finite or infinite. We will deal in most of the paper with the infinite horizon case but all of the results carry over immediately to finitely repeated games if  $T$  is large enough. In the repeated game the overall payoff for player  $i$  from period  $t$  onwards (and including period  $t$ ) is given by

$$(1) \quad V_i^t = \sum_{\tau=t}^{\infty} \delta_i^{\tau-t} g_i^\tau,$$

where  $\delta_i$  denotes her (his) discount factor ( $0 \leq \delta_i < 1$ ). Our results are stated in terms of average discounted payoffs  $v_i$ , where

$$(2) \quad v_i = (1 - \delta_i) \cdot V_i^0 = (1 - \delta_i) \cdot \sum_{\tau=0}^{\infty} \delta_i^\tau g_i^\tau.$$

After each period both players observe the actions that have been taken. They have perfect recall and can condition their play on the entire past history of the game. Let  $h^t$  be a specific history of the repeated game out of the set  $H^t = (A_1 \times A_2)^t$  of all possible histories up to and including period  $t$ . A pure strategy  $s_i$  for player  $i$  in the repeated game is a sequence of maps  $s_i^t : H^{t-1} \rightarrow A_i$ . Correspondingly, let  $\sigma_i = (\sigma_i^1, \sigma_i^2, \dots)$  denote a mixed (behavioral) strategy of player  $i$ , where  $\sigma_i^t : H^{t-1} \rightarrow \mathcal{A}_i$ . For notational convenience the dependence on history is suppressed if there is no ambiguity. The set of all pure (mixed) strategies is denoted by  $S_i$  ( $\Sigma_i$  respectively).

Let  $B : \mathcal{A}_1 \mapsto \mathcal{A}_2$  be the best response correspondence of player two in the stage game and define

$$(3) \quad g_1^* = \max_{a_1 \in A_1} \min_{\alpha_2 \in B(a_1)} g_1(a_1, \alpha_2)$$

as the ‘‘commitment payoff’’ of player one. That is  $g_1^*$  is the most player one could guarantee for herself in the stage game if she could commit to any pure strategy  $a_1 \in A_1$ . Note that the minimum over all  $\alpha_2 \in B(a_1)$  has to be taken since player two may be indifferent between several best responses to  $a_1$  in which case he may take the response

player one prefers least.<sup>3</sup> Let  $a_1^*$  (her “commitment action”) satisfy

$$(4) \quad \min_{\alpha_2 \in B(a_1^*)} g_1(a_1^*, \alpha_2) = g_1^* .$$

Furthermore, let  $\alpha_2^* \in B(a_1^*)$  denote any strategy of player two which is a best response to  $a_1^*$  and define

$$(5) \quad g_2^* = g_2(a_1^*, \alpha_2^*) .$$

So  $g_2^*$  is the most player two would get in the stage game if player one were committed to  $a_1^*$ . Suppose  $B(a_1^*) \neq \mathcal{A}_2$  (otherwise the game is “trivial” because player one’s commitment payoff is her maxmin payoff). Then there exists a  $\tilde{a}_2 \notin B(a_1^*)$  such that

$$(6) \quad \tilde{g}_2 = g_2(a_1^*, \tilde{a}_2) = \max_{a_2 \notin B(a_1^*)} g_2(a_1^*, a_2) < g_2^* .$$

Note that the maximum exists because it is taken over the finite set of all (pure) actions  $a_2 \notin B(a_1^*)$ . So  $\tilde{g}_2$  is the maximum player two can get if he does not take an action which is a best response against  $a_1^*$ , given that player one takes her commitment action. Finally, define the maximal payoff player two can get at all as

$$(7) \quad \bar{g}_2 = \max_{a_2 \in \mathcal{A}_2} \max_{a_1 \in \mathcal{A}_1} g_2(a_1, a_2) .$$

Clearly, in the repeated game it must be true that

$$(8) \quad V_2^t \leq \sum_{\tau=t}^{\infty} \delta_2^{\tau-t} \cdot \bar{g}_2 = \frac{\bar{g}_2}{1 - \delta_2} = \bar{V}_2^t$$

for all  $t$  and all  $h^{t-1} \in H^{t-1}$ .

Consider now a perturbation of this complete information game such that in period 0 (before the first stage game is played) the “type” of player one is drawn by nature out of a countable set  $\Omega = (\omega_0, \omega_1, \dots)$  according to the probability measure  $\mu$ . Player one’s payoff function now additionally depends on her type, so  $g_1 : A_1 \times A_2 \times \Omega \rightarrow \mathbb{R}$ . The perturbed game  $G^T(\mu)$  is a game with incomplete information in the sense of Harsanyi (1967-68). In the perturbed game a strategy of player one may not only depend on history

---

<sup>3</sup>Fudenberg and Levine (1989) refer to  $g_1^*$  as the “Stackelberg payoff”. However, it is now customary to use this expression only for  $\max_{a_1} \max_{\alpha_2 \in B(a_1)} g_1(a_1, \alpha_2)$ , that is for the maximum payoff player one could get if he could publicly commit himself to any action  $a_1$  and player two chooses the best response player one prefers *most*. See Fudenberg (1990). The analysis can be extended to the more general case where player one would like to commit himself to a mixed strategy or to a strategy dependent on history. See Fudenberg and Levine (1992) and the remarks in Section 6.

but also on her type, so  $\sigma_1^t : H^{t-1} \times \Omega \rightarrow \mathcal{A}_1$ . Two types out of the set  $\Omega$  are of particular importance:

- The “normal” type of player one is denoted by  $\omega_0$ . Her payoff function is the same as in the unperturbed game:

$$(9) \quad g_1(a_1, a_2, \omega_0) = g_1(a_1, a_2).$$

In many applications  $\mu(\omega_0)$  will be close to 1. However, we have to require only that  $\mu(\omega_0) = \mu^0 > 0$ .

- The “commitment” type is denoted by  $\omega^*$ . For her it is a dominant strategy in the repeated game always to play  $a_1^*$ . This is for example the case if her payoff function satisfies

$$(10) \quad g_1(a_1^*, a_2, \omega^*) = g_1(a_1^*, a_2', \omega^*) > g_1(a_1, a_2', \omega^*)$$

for all  $a_1 \neq a_1^*$ ,  $a_1 \in A_1$ , and all  $a_2, a_2' \in A_2$ . The dominant strategy property in the repeated game implies that in any Nash equilibrium player one with type  $\omega^*$  has to play  $a_1^*$  in every period along the equilibrium path. This in turn implies that if  $\mu(\omega^*) = \mu^* > 0$  then with positive probability there exists a history in any Nash equilibrium with  $s_1^t = a_1^*$  for all  $t$ . The set of all such histories is denoted by  $H^*$ .

We will now restate an important lemma of Fudenberg-Levine (1989) about statistical inference which is basic to the following analysis. The lemma says that if  $\omega^*$  has positive probability and if player two observes  $a_1^*$  being played in every period then there is a fixed finite upper bound on the number of periods in which player two will believe  $a_1^*$  is “unlikely” to be played. The intuition for this result is the following. Consider any history  $h^{t-1} \in H^*$  in which player one has always played  $a_1^*$  up to period  $t - 1$ . Suppose player two believes that the probability of  $a_1^*$  being played in period  $t$  is smaller than  $\bar{\pi}$ ,  $0 \leq \bar{\pi} < 1$ . If player two observes  $a_1^*$  being played in  $t$  he is “surprised” to some extent and will update his beliefs. Because the commitment type chooses  $a_1^*$  with probability 1 while player two expected  $a_1^*$  to be played with a probability bounded away from 1 it follows from Bayes’ law that the updated probability that he faces the commitment type has to increase by an amount bounded away from 0. However, this cannot happen arbitrarily often because the updated probability of the commitment type cannot become



bigger than 1. This gives the upper bound on the number of periods in which player two may expect  $a_1^*$  to be played with a probability less than  $\bar{\pi}$ . Note that this argument is independent of the discount factors of the two players.

To put it more formally: Each (possibly mixed) strategy profile  $(\sigma_1, \sigma_2)$  induces a probability distribution  $\pi$  over  $(A_1 \times A_2)^\infty \times \Omega$ . Given a history  $h^{t-1}$  let  $\pi^t(a_1^*)$  be the probability attached by player two to the event that the commitment strategy is being played in period  $t$ , i.e.  $\pi^t(a_1^*) = \text{Prob}(s_1^t = a_1^* \mid h^{t-1})$ . Note that since  $h^{t-1}$  is a random variable  $\pi^t(a_1^*)$  is a random variable as well. Fix any  $\bar{\pi}$ ,  $0 \leq \bar{\pi} \leq 1$ , and consider any history  $h$  induced by  $(\sigma_1, \sigma_2)$ . Along this history let  $n(\pi^t(a_1^*) \leq \bar{\pi})$  be the number (possibly infinite) of the random variables  $\pi^t(a_1^*)$  for which  $\pi^t(a_1^*) \leq \bar{\pi}$ . Again, since  $h$  is a random variable, so is  $n$ .

**Lemma 1** *Let  $0 \leq \bar{\pi} < 1$ . Suppose  $\mu(\omega^*) = \mu^* > 0$ , and that  $(\sigma_1, \sigma_2)$  are such that  $\text{Prob}(h \in H^* \mid \omega^*) = 1$ . Then*

$$(11) \quad \text{Prob} \left[ n(\pi^t(a_1^*) \leq \bar{\pi}) > \frac{\log \mu^*}{\log \bar{\pi}} \mid h \in H^* \right] = 0.$$

*Furthermore, for any infinite history  $h$  such that the truncated histories  $h_t$  all have positive probability and such that  $a_1^*$  is always played,  $\mu(\omega^* \mid h_t)$  is nondecreasing in  $t$ .*

Proof: See Fudenberg-Levine (1989), Lemma 1.

One feasible strategy for player one with type  $\omega_0$  is of course to mimic the commitment type and always to play  $a_1^*$ . Lemma 1 does not say that in this case  $\mu(\omega^* \mid h_t \in h^*)$  converges to 1, i.e. that player two will gradually become convinced that he is facing  $\omega^*$  if he observes  $a_1^*$  always being played. Rather it says that if he observes  $a_1^*$  being played in every period he cannot continue to believe that  $a_1^*$  is “unlikely” to be played.

Suppose that player two is completely myopic, that is he is only interested in his payoff of the current period. Fudenberg and Levine show that there is a  $\bar{\pi} < 1$  such that if the probability that player one will play  $a_1^*$  is bigger than  $\bar{\pi}$  then a short-run player two will choose a best response against  $a_1^*$ . Thus, if player one mimics the commitment type, then by Lemma 1 her short-run opponents will take  $a_2 \notin B(a_1^*)$  in at most  $k = \frac{\log \mu^*}{\log \bar{\pi}}$

periods. The worst that can happen to player one is that these  $k$  periods occur in the beginning of the game and that in each of these periods she gets

$$(12) \quad \underline{g}_1 = \min_{\alpha_2 \in \mathcal{A}_2} g_1(a_1^*, \alpha_2).$$

This argument provides the intuition for the following theorem.

**Theorem 1 (Fudenberg-Levine)** *Let  $\delta_2 = 0$ ,  $\mu(\omega^0) > 0$ , and  $\mu(\omega^*) = \mu^* > 0$ . Then there is a constant  $k(\mu^*)$  otherwise independent of  $(\Omega, \mu)$ , such that*

$$(13) \quad v_1(\delta_1, \mu^*; \omega^0) \geq (1 - \delta_1^{k(\mu^*)}) \cdot \underline{g}_1 + \delta_1^{k(\mu^*)} \cdot g_1^*,$$

where  $v_1(\delta_1, \mu^*; \omega^0)$  is any average equilibrium payoff of player one with type  $\omega_0$  in any Nash equilibrium of  $G^\infty(\mu)$ .

If  $\delta_1$  goes to 1 the “normal” type of player one can guarantee herself on average at least her commitment payoff no matter what other types may be around with positive probability. The result is discussed in more detail in Fudenberg and Levine (1989). Note however that Theorem 1 is crucially based on the assumption that player two is completely myopic. If he cares about future payoffs then he may trade off short-run losses against long-run gains. Thus, even if he believes that  $a_1^*$  will be played with a probability arbitrarily close or equal to 1, he may take an action  $a_2$  which is not a short-run best response against  $a_1^*$ . One intuitive reason for this could be that he might invest in screening the different types of player one. Even if this yields losses in the beginning of the game the investment may well pay off in the future. This leads Fudenberg and Levine to conclude that their result does not apply to two long-run player games. The main point of our paper, however, is to show that for a more restricted class of games a similar result holds in the two long-run players case as well. Since player two’s discount factor is smaller than 1 the returns from an investment may not be delayed too far to the future. He will not “test” player one’s type arbitrarily often if the probability that she will play  $a_1^*$  is always arbitrarily close to one. This idea will be used in Section 4 to prove an analog of Theorem 1 for two long-run player games.

### 3 A Game not of Conflicting Interests

Before establishing our main result let us show that Theorem 1 does not carry over to all repeated games with two long-run players. We give a counterexample of a game in which the normal type of player one cannot guarantee herself almost her commitment payoff in all Nash equilibria. The example is instructive for two reasons. First, it shows how to construct an equilibrium in which the normal type of player one gets strictly less than her commitment payoff. This equilibrium is not only a Nash but a sequential equilibrium which survives all standard refinements. Second, the construction leads to a necessary and sufficient condition on the class of games for which Theorem 1 can be generalized to the two long-run player case.

Consider an infinite repetition of the following stage game with three types of player one:

	<i>L</i>	<i>R</i>		<i>L</i>	<i>R</i>		<i>L</i>	<i>R</i>		
<i>U</i>	10	10	0	10	0	<i>U</i>	1	10	1	0
<i>D</i>	0	0	1	1	1	<i>D</i>	1	0	1	1
	“normal” type			“commitment” type			“indifferent” type			
	$\mu^0 = 0.8$			$\mu^* = 0.1$			$\mu^i = 0.1$			

FIGURE 1 – A game with common interests.

Player one chooses between  $U$  and  $D$  and her payoff is given in the upper left corners of each cell. Clearly the normal type of player one would like to publicly commit always to play  $U$  which would give her a commitment payoff of 10 per period in every Nash equilibrium. For the commitment type it is indeed a dominant strategy in the repeated game always to play  $U$ . The indifferent type, however, is indifferent between  $U$  and  $D$  no matter what player two does.

If  $0.75 \leq \delta_1 \leq 1$  and  $0.95 \leq \delta_2 \leq 1$  then the following strategies and beliefs form a

sequential equilibrium of  $G^\infty$  which gives the normal type of player one

$$(14) \quad \lim_{\delta_1 \rightarrow 1} v_1(\omega^0) = 9.5 < 10 = g_1^* .$$

- *Normal type of player one:* “Play  $U$ . If you ever played  $D$ , switch to playing  $D$  forever.”
- *Commitment type of player one:* “Always play  $U$ .”
- *Indifferent type of player one:* “Always play  $U$  along the equilibrium path. If there has been any deviation by any player in the past switch to playing  $D$  forever.”
- *Player two:* “Alternate between 19 times  $L$  and 1 times  $R$  along the equilibrium path. If player one ever played  $D$ , switch to  $R$  forever. If player two himself deviated in the last period, play  $L$  in the following period. If player one reacted to the deviation by playing  $U$ , go on playing  $L$  forever. If she reacted with  $D$ , play  $R$  forever.”
- *Beliefs:* Along the equilibrium path beliefs don't change. If player two ever observes  $D$  to be played he puts probability 0 on the commitment type. If player one reacts to a deviation of player two by playing  $U$  the indifferent type gets probability 0. In both cases the respective two other types may get arbitrary probabilities which add up to 1.

Why is this an equilibrium? Consider the normal type of player one. Clearly she would like to signal that she is the normal or the commitment type. Since all three types of player one always play  $U$  along the equilibrium path the only way to transmit information about her type is to play  $D$ . However, playing  $D$  “kills” the commitment type, because for her it is a dominant strategy always to play  $U$ . But without the commitment type it is impossible to get rid of the “bad” equilibrium  $(D, R)$ . What about player two? He expects  $U$  always to be played along the equilibrium path. Nevertheless he plays  $R$ , which is not a short-run best response, in every twentieth period. His problem is that he faces the indifferent type with positive probability. If he chooses  $L$  when he is supposed to play  $R$ , then this might trigger a continuation equilibrium against the indifferent type

which gives him far less than what he would have got from playing against the normal or commitment type of player one. It is this risk which sustains the equilibrium outcome.

Note that there are very few restrictions imposed on the updating of beliefs in information sets which are not reached on the equilibrium path. The example only requires that if  $D$  is played for the first time the commitment type gets probability 0, which is perfectly reasonable given that for her it is a dominant strategy in the repeated game always to play  $U$ .

To what extent does the example rely on the existence of the indifferent type? Without the indifferent type it is still possible to construct a Nash equilibrium which gives player one less than her commitment payoff. Actually, this is very simple: The normal and the commitment type of player one always play  $U$  along the equilibrium path. After any deviation they switch to playing  $D$  forever. Player two alternates playing one period  $L$  and one period  $R$ . If there has been any deviation, he plays  $R$  forever. Note that the average payoff of the normal type of player one is only 5. This clearly is a Nash equilibrium, but it is not sequential. It requires for example that the commitment type plays  $D$  off the equilibrium path.<sup>4</sup>

What sustains both equilibria is the possibility of a continuation equilibrium which punishes player two if he plays his short-run best response against  $a_1^*$  in periods when he is supposed not to do so. Note that this construction does not work if player two is already hold down to his minimax payoff by the commitment strategy of player one, since in this case nothing worse can happen to him. In the next section we show that this is the only case in which Fudenberg and Levine's result can be generalized to the two long-run player case.

---

<sup>4</sup>Whether there exists a *sequential* equilibrium in which player one gets substantially less than 10 if there are only the normal and the commitment type around is an open question. Note, however, that we want to characterize equilibrium outcomes which are robust to general perturbations of the informational structure of the game. From this perspective it makes little sense to restrict attention to two possible types only.

## 4 Main Results

### 4.1 The Theorem

Suppose that the commitment strategy of player one holds player two down to his minimax payoff. In this case there is no “risk” in playing a best response against  $a_1^*$  because player two cannot get less than his minimax payoff in any continuation equilibrium. This motivates the following definition:

**Definition 1** *A game  $g$  is called a game of conflicting interests with respect to player one if the commitment strategy of player one holds player two down to his minimax payoff, i.e. if*

$$(15) \quad g_2^* = g_2(a_1^*, \alpha_2^*) = \min_{\alpha_1} \max_{\alpha_2} g_2(\alpha_1, \alpha_2).$$

“Conflicting interests” are a necessary and sufficient condition for our main result. Note that the definition puts no restriction on the possible perturbations of the payoffs of player one. It is a restriction only on the commitment strategy and on the payoff function of player two. We will discuss this class of games extensively and give several examples in Section 5. Clearly, in a game with conflicting interests player two can guarantee himself in any continuation equilibrium after any history  $h_t$  at least

$$(16) \quad V_2^* = \frac{1}{1 - \delta_2} \cdot g_2^*.$$

This is crucial to establish the following result:

**Lemma 2** *Let  $g$  be a game of conflicting interests with respect to player one and let  $\mu(\omega^*) = \mu^* > 0$ . Consider any Nash equilibrium  $(\hat{\sigma}_1, \hat{\sigma}_2)$  and any history  $h^t$  consistent with this equilibrium in which player one has always played  $a_1^*$ . Suppose that, given this history, the equilibrium strategy of player two prescribes to take  $s_2^{t+1} \notin B(a_1^*)$  with positive probability in period  $t + 1$ . For any  $\delta_2$ ,  $0 < \delta_2 < 1$ , there exists a finite integer  $M$ ,*

$$(17) \quad M \geq N = \frac{\ln(1 - \delta_2) + \ln(g_2^* - \tilde{g}_2) - \ln(\bar{g}_2 - \tilde{g}_2)}{\ln \delta_2} > 0,$$

and a positive number  $\epsilon$ ,

$$(18) \quad \epsilon = \frac{(1 - \delta_2)^2 \cdot (g_2^* - \tilde{g}_2)}{\bar{g}_2 - \tilde{g}_2} - \delta_2^M \cdot (1 - \delta_2) > 0 ,$$

such that in at least one of the periods  $t + 1, t + 2, \dots, t + M$  the probability that player one does not take  $a_1^*$ , given that he always played  $a_1^*$  before, must be at least  $\epsilon$ .

Proof: See appendix.

Let us briefly outline the intuition behind this result. Because  $g$  is of conflicting interests player two can guarantee himself at least  $\underline{V}_2^{t+1} = \frac{g_2^*}{1-\delta_2}$  in any continuation equilibrium after any history  $h_t$ . Therefore, if he tries to test player one's type and takes an action  $s_2^{t+1} \notin B(a_1^*)$  in period  $t + 1$  this must give him an expected payoff of at least  $\underline{V}_2^{t+1}$  for the rest of the game. If player one chooses  $a_1^*$  with a probability arbitrarily close or equal to 1, then choosing an action  $a_2 \notin B(a_1^*)$  yields a "loss" of at least  $g_2^* - \tilde{g}_2 > 0$  in this period. Recall that  $\tilde{g}_2$  is defined as the maximal payoff player two gets if he does not take a best response against  $a_1^*$ . On the other hand,  $\bar{g}_2$  is an upper bound on what player two may get in any period in which player one does not take her commitment action, and - of course - he cannot get more than  $g_2^*$  if she plays  $a_1^*$ . But if future payoffs are bounded and  $\delta_2 < 1$  then the compensation for an expected loss today must not be delayed too far to the future. The numbers  $M$  and  $\epsilon$  are constructed such that if player two takes a strategy  $s_2^{t+1} \notin B(a_1^*)$  in period  $t + 1$  then it cannot be true that in each of the next  $M$  periods the probability that player one takes her commitment action is bigger than  $(1 - \epsilon)$ . Otherwise player two would get less than his minimax payoff in equilibrium, a contradiction. Note that this argument also holds for finitely repeated games if  $T$  is large enough.

Lemma 2 holds in any proper subform of  $G$  as long as player one always played  $a_1^*$  in the history up to that subform. Thus if player two chooses actions  $a_2 \notin B(a_1^*)$  along  $h^*$  in  $n \cdot M$  periods, then in at least  $n$  of these periods the probability that player one does not play  $a_1^*$  must be at least  $\epsilon$ . Together with Lemma 1 this implies our main theorem:

**Theorem 2** *Let  $g$  be of conflicting interests with respect to player one and let  $\mu(\omega^0) > 0$ , and  $\mu(\omega^*) = \mu^* > 0$ . Then there is a constant  $k(\mu^*, \delta_2)$  otherwise*

independent of  $(\Omega, \mu)$ , such that

$$(19) \quad v_1(\delta_1, \delta_2, \mu^*; \omega^0) \geq \left(1 - \delta_1^{k(\mu^*, \delta_2)}\right) \cdot \underline{g}_1 + \delta_1^{k(\mu^*, \delta_2)} \cdot g_1^*,$$

where  $v_1(\delta_1, \delta_2, \mu^*; \omega^0)$  is any average equilibrium payoff of player one with type  $\omega_0$  in any Nash equilibrium of  $G^\infty(\mu)$ .

Proof: Consider the strategy for the normal type of player one of always playing  $a_1^*$ . Take the integer  $M = [N] + 1$ , where  $[N]$  is the integer part of  $N$ , and a real number  $\epsilon > 0$ , where  $N$  and  $\epsilon$  are defined in Lemma 2. By Lemma 2 we know that if player two takes an action  $a_2 \notin B(a_1^*)$  then there is at least one period (call it  $\tau_1$ ) among the next  $M$  periods in which the probability that player one will play  $a_1^*$  (denoted by  $\pi_{\tau_1}^*$ ) is smaller than  $(1 - \epsilon)$ . So

$$(20) \quad \pi_{\tau_1}^* < 1 - \epsilon \equiv \bar{\pi}.$$

However, by Lemma 1 we know that

$$(21) \quad \pi \left[ n(\pi_t^* \leq \bar{\pi}) > \frac{\ln \mu^*}{\ln \bar{\pi}} \mid h^* \right] = 0.$$

That is, the probability that player one takes her commitment action cannot be smaller than  $\bar{\pi}$  in more than  $\frac{\ln \mu^*}{\ln \bar{\pi}}$  periods. Therefore, player two cannot choose actions  $a_2 \notin B(a_1^*)$  more often than

$$(22) \quad k = M \cdot \frac{\ln \mu^*}{\ln(1 - \epsilon)}$$

times. Substituting  $M = [N] + 1$  and  $\epsilon$  from Lemma 2, we get

$$(23) \quad k(\mu^*, \delta_2) = ([N] + 1) \cdot \frac{\ln \mu^*}{\ln \left( 1 - \frac{(1 - \delta_2) \cdot (g_2^* - \underline{g}_2)}{\bar{g}_2 - \underline{g}_2} + \delta_2^{[N] + 1} \right)}.$$

In the worst case player two chooses these actions in the first  $k(\mu^*, \delta_2)$  periods. This gives the lower bound of the theorem. Q.E.D.

If  $\delta_1 \rightarrow 1$  (keeping  $\delta_2$  fixed) then the equilibrium payoff of the normal type of player one is bounded below by her commitment payoff. Thus, in the limit our theorem gives the same lower bound as Fudenberg and Levine's theorem does for the case of a long-run player facing a sequence of short-run opponents. Their result can be obtained as a special



case of Theorem 2 for the class of games with conflicting interests. Note that if  $\delta_2$  goes to 0 then  $N$  goes to 0. So

$$(24) \quad \lim_{\delta_2 \rightarrow 0} k(\mu^*, \delta_2) = \frac{\ln \mu^*}{\ln \left(1 - \frac{g_2^* - \tilde{g}_2}{\bar{g}_2 - \tilde{g}_2}\right)} = \frac{\ln \mu^*}{\ln \left(\frac{\bar{g}_2 - g_2^*}{\bar{g}_2 - \tilde{g}_2}\right)}.$$

In a game with conflicting interests a short-run player two will play a best response against  $a_1^*$  if

$$(25) \quad g_2^* > \pi \cdot \tilde{g}_2 + (1 - \pi) \cdot \bar{g}_2$$

or, equivalently, if

$$(26) \quad \pi > \frac{\bar{g}_2 - g_2^*}{\bar{g}_2 - \tilde{g}_2} \equiv \bar{\pi}.$$

Using (26) in Lemma 1 immediately implies Theorem 1.

It is important to note that the lower bound given in Theorem 2 depends on the discount factor of player 2. If  $\delta_2$  increases, so does  $k(\mu^*, \delta_2)$ , and the lower bound is reduced. Hence, to obtain his commitment payoff in equilibrium player one has to be sufficiently patient as compared to player two. The following corollaries elaborate on the importance of the relative patience of the players.

**Corollary 1** *For any  $\delta_2 < 1$ ,  $\mu^* > 0$  and  $\epsilon > 0$  there exists a  $\underline{\delta}_1(\delta_2, \mu^*, \epsilon) < 1$ , such that for any  $\delta_1 \geq \underline{\delta}_1(\delta_2, \mu^*, \epsilon)$  the average payoff of the normal type of player one is at least  $g_1^* - \epsilon$ .*

Proof: Choose  $\underline{\delta}_1$  such that

$$(27) \quad g_1^* - \epsilon = \left(1 - \underline{\delta}_1^{k(\mu^*, \delta_2)}\right) \cdot \underline{g}_1 + \underline{\delta}_1^{k(\mu^*, \delta_2)} \cdot g_1^*.$$

Solving for  $\underline{\delta}_1$  yields

$$(28) \quad \underline{\delta}_1 = \underline{\delta}_1(\mu^*, \delta_2, \epsilon) = \sqrt[k(\mu^*, \delta_2)]{\frac{g_1^* - \underline{g}_1 - \epsilon}{g_1^* - \underline{g}_1}} < 1.$$

Clearly, if  $\delta_1 \geq \underline{\delta}_1(\mu^*, \delta_2, \epsilon)$  then

$$(29) \quad \begin{aligned} v_1(\delta_1, \delta_2, \mu^*; \omega^0) &\geq \left(1 - \delta_1^{k(\mu^*, \delta_2)}\right) \cdot \underline{g}_1 + \delta_1^{k(\mu^*, \delta_2)} \cdot g_1^* \\ &\geq \left(1 - \underline{\delta}_1^{k(\mu^*, \delta_2)}\right) \cdot \underline{g}_1 + \underline{\delta}_1^{k(\mu^*, \delta_2)} \cdot g_1^* = g_1^* - \epsilon. \end{aligned}$$

*Q.E.D.*

**Corollary 2** Consider any sequence  $\{\delta_2^n\}$ ,  $\delta_2^n < 1$ ,  $\lim_{n \rightarrow \infty} \delta_2^n = 1$  and fix  $\epsilon > 0$ . Then there exists a sequence  $\{\delta_1^n\}$ ,  $\delta_1^n < 1$ ,  $\lim_{n \rightarrow \infty} \delta_1^n = 1$ , such that for any  $\{\delta_1^n, \delta_2^n\}$  the average payoff of the normal type of player one is bounded below by  $g_1^* - \epsilon$ .

Proof: Take any  $\{\delta_2^n\}$  and fix  $\epsilon > 0$ . Choose  $\{\delta_1^n\} \rightarrow 1$  such that  $\delta_1^n > \underline{\delta}_1(\delta_2, \mu^*, \epsilon)$  for all  $n$ , where  $\underline{\delta}_1(\delta_2, \mu^*, \epsilon)$  is given by (28). Then the result follows immediately from the previous corollary. Q.E.D.

Corollary 2 shows that there is an area in the  $\delta_1 - \delta_2$  space such that for any sequence  $\{\delta_1^n, \delta_2^n\} \rightarrow (1, 1)$  in this area player one gets at least her commitment payoff (up to an arbitrarily small  $\epsilon$ ) for any pair of discount factors along this sequence. Note, however, that  $\lim_{n \rightarrow \infty} \frac{1-\delta_1}{1-\delta_2} = 0$ , i.e. in the limit player one is infinitely more patient than player two. This observation helps to understand a related result of Cripps and Thomas (1991) who consider repeated games without discounting, in which players maximize the limit of the mean of their payoffs. Under slightly stronger conditions on the possible perturbations they show that if the game has conflicting interests and if there is a positive prior probability of a commitment type, then player one gets at least her commitment payoff as the Banach limit of the mean of her stage game payoffs. However, the case of no discounting obscures the role of the relative patience of the players. We can give examples of equilibria in games with conflicting interests where  $\delta_1 \rightarrow 1$ ,  $\delta_2 \rightarrow 1$ ,  $\lim_{n \rightarrow \infty} \frac{1-\delta_1}{1-\delta_2} > 0$ , and player one's equilibrium payoff is bounded away from her commitment payoff for any  $\{\delta_1^n, \delta_2^n\}$  along this sequence. Thus, if player one is not patient enough as compared to player two our lower bound does not apply.

If player one has to be much more patient than player two the reader might be left with the impression that we are essentially back to Fudenberg and Levine (1989) where a long-run player faces a sequence of short-run players. However, this is not the case. First, Fudenberg and Levine's result requires  $\delta_2 = 0$  while here  $\delta_2$  may be arbitrarily close to 1. Second, we are going to show in the next subsection that whenever the game is not of conflicting interests it is possible to find an equilibrium which violates Fudenberg and Levine's lower bound no matter how much more patient player one is as compared to player two. Thus, there is a fundamental difference between repeated games in which one player does not care at all about her future payoffs and games in which she does care

but is less patient than her opponent. Finally, the importance of the relative patience of the two players is very intuitive as will be shown after we have introduced the case of two-sided uncertainty in subsection 4.3.

## 4.2 Necessity of the “Conflicting Interests” Condition

The question arises whether Theorem 2 also holds for games which are not of conflicting interests. If the game is not “trivial” in the sense that player one’s commitment payoff is equal to her minimax payoff<sup>5</sup> the answer is no:

**Proposition 1** *Let  $g$  be a non-trivial game which is not of conflicting interests. Then for any  $\epsilon > 0$  there is an  $\eta > 0$  and a  $\underline{\delta}_2 < 1$  such that the following holds: There is a perturbation of  $g$ , in which the commitment type of player 1 has positive probability and the normal type has probability  $(1 - \epsilon)$ , and there is a sequential equilibrium of this perturbed game, such that the limit of the average payoff of the normal type of player one for  $\delta_1 \rightarrow 1$  is bounded away from her commitment payoff by at least  $\eta$  for any  $\delta_2 > \underline{\delta}_2$ .*

Proof: See appendix.

Proposition 1 shows that the condition of conflicting interests is not only sufficient but also necessary for Theorem 2 to hold, in fact, it is a little bit stronger than that in two respects. First, it says that if the game is not of conflicting interests, then it is not only possible to find a Nash equilibrium which violates Fudenberg and Levine’s lower bound, but also to find a sequential equilibrium. As has been indicated in Section 3, the construction of a Nash equilibrium using threats which are not credible is much simpler. Secondly, Theorem 2 only requires that  $\mu(\omega^0) > 0$  in the perturbed game. So we could have established necessity by constructing a perturbation which gives a high prior probability to an “indifferent” type who credibly threatens to punish any deviation of player two from the equilibrium path we want to sustain. However, in many economic applications it is natural to assume that  $\mu(\omega^0)$  is close to one. This is why we provide a

---

<sup>5</sup>It is well known that a player can always guarantee herself at least her minimax payoff in any Nash-equilibrium.

stronger proposition which says that even if  $\mu(\omega^0)$  is arbitrarily close to one it is possible to construct a sequential equilibrium in which the payoff of the normal type of player one is bounded away from  $g_1^*$ .

Note that in Proposition 1  $\delta_1 \rightarrow 1$  while  $\delta_2$  is fixed, so player one may be arbitrarily more patient than player two. Thus, Proposition 1 shows that there is an important difference between games with two long run players, one of whom is more patient than the other, and games in which a long-run player faces a sequence of short-run players. In the latter Fudenberg and Levine's bound holds for for all stage games, in the former it only holds for games with conflicting interests.

### 4.3 Two-sided Incomplete Information and Two-sided Conflicting Interests

If there are two long-run players it is most natural to ask what happens if there is two-sided uncertainty. Our result can be extended to this case in the following way. Suppose the game is perturbed such that there is incomplete information about both the payoff functions of player one and player two. Let  $\omega_i$  denote player  $i$ 's type which is drawn by nature in the beginning of the game out of the countable set  $\Omega_i$  according to the probability measure  $\mu_i$ ,  $i \in \{1, 2\}$ . Let  $\omega_i^0$  and  $\omega_i^*$  represent the normal and the commitment types, respectively. Finally, suppose that the game is of conflicting interests with respect to player  $i$ , i.e. player  $i$ 's commitment strategy holds player  $j$  down to his minimax payoff. Without loss of generality let  $i = 1$ . Now consider the normal type of player two. In the proof of Lemma 2 we didn't say why player two might choose an action which is not a best response against player one's commitment strategy. He might do so because he wants to test player one's type or because he wants to build up a reputation himself. No matter what the reason is, Lemma 2 states that if he takes  $a_2 \notin B(a_1^*)$ , then he must expect that player one chooses  $s_1^t \neq a_1^*$  in one of the following periods with strictly positive probability. This holds for the normal type of player two no matter what other possible types of player two exist with positive probability.

A possible strategy of player one still is to play  $a_1^*$  in every period. If she faces the normal type of player two, then by Theorem 2 there are at most  $k(\mu_1^*, \delta_2)$  periods in which player two will not play a best response against  $a_1^*$ . In the worst case for player

one this happens in the first  $k$  periods of the game. On the other hand, if she does not face the normal type of player two her expected payoff is at least  $\underline{g}_1$  in every period. This argument establishes a lower bound for the expected payoff of the normal type of player  $i$  which is given in the following theorem:

**Theorem 3** *Let  $g$  be of conflicting interests with respect to player  $i$  and let  $\mu_i(\omega_i^0) = \mu_i^0 > 0$  and  $\mu_i(\omega_i^*) = \mu_i^* > 0$ ,  $i \in \{1, 2\}$ . Then there are constants  $k_i(\mu_i^*, \delta_j)$  otherwise independent of  $(\Omega_i, \Omega_j, \mu_i)$ , such that*

$$(30) \quad v_i(\delta_i, \delta_j, \mu_i^*, \mu_j^0; \omega_i^0) \geq \left(1 - \mu_j^0 \delta_i^{k_i(\mu_i^*, \delta_j)}\right) \underline{g}_i + \mu_j^0 \delta_i^{k_i(\mu_i^*, \delta_j)} g_i^*$$

where  $v_i(\delta_1, \delta_2, \mu_i^*, \mu_j^0; \omega_i^0)$  is any average equilibrium payoff of player  $i$  with type  $\omega_i^0$  in any Nash equilibrium of  $G^\infty(\mu)$ .

Thus, if the probability of the normal type of player two is close to 1 and if player one is very patient, then the lower bound for her average payoff is again close to her commitment payoff.

What can be said if  $g$  has two-sided conflicting interests, i.e. if each player would like to commit to a strategy which holds his opponent down to his minimax payoff. Of course, if there are two-sided conflicting interests and if both players are equally patient it is impossible that each of them gets his most preferred payoff. But suppose that  $\delta_i$  and  $\delta_j$  differ. The bigger player  $j$ 's discount factor the bigger is  $k_i(\mu_i^*, \delta_j)$ , i.e. the number of periods in which player  $i$  must expect that a strategy other than the best response against her commitment strategy is played, and the lower is her lower bound. On the other hand, if  $\delta_j$  is kept fixed and  $\delta_i$  goes to 1 then this  $k$  periods become less and less important, and in the limit player  $i$  will get his commitment payoff. This is very intuitive. In a symmetric game with conflicting interests reputation building has an effect only if one of the parties is sufficiently more patient than the other.

Theorems 2 and 3 are in striking contrast to the message of the Folk theorem for games with incomplete information by Fudenberg and Maskin (1986). The Folk theorem says that any feasible payoff vector which gives each of the players at least his minimax payoff can be sustained as an equilibrium outcome of the perturbed game if the right

perturbation is chosen.<sup>6</sup> Theorems 2 and 3 show that this result is not robust against further perturbations. If one of the players is patient enough and if her commitment type has positive probability then - no matter what other types are around with positive probability - Theorem 2 imposes a tight restriction on the set of equilibrium outcomes in any Nash equilibrium.<sup>7</sup>

We have to be very precise here in what is meant by robustness. Fudenberg (1990) argues that the Folk theorem is robust against a further perturbation of the informational structure in the following sense: Consider a sequential equilibrium which has been constructed using an  $\epsilon$  probability of a “crazy” type as the Folk theorem suggests. If for a given time horizon of the game other types are introduced with a probability which is small as compared to  $\epsilon$ , then this is still an equilibrium. However, if the time horizon goes to infinity the equilibrium must break down. We have shown, that if the game is of conflicting interests and if there is an arbitrarily small probability of a commitment type, then the commitment type will dominate the play as  $T$  becomes large enough. Thus, if we are interested in the set of equilibria for  $T \rightarrow \infty$ , the Folk theorem is not robust against small perturbations of the informational structure.

## 5 Examples

### 5.1 The Chain Store Game

Consider the classical chain store game, introduced by Selten (1978), with two long-run players. In every period the entrant may choose to enter a market ( $I$ ) or to stay out ( $O$ ),

---

<sup>6</sup>Fudenberg and Maskin’s Folk theorem for games with incomplete information considers only finitely repeated games without discounting. However, the extension to discounting and an infinite horizon is straightforward.

<sup>7</sup>Note that even if the game is not of conflicting interests we can still impose some restriction on the set of Nash equilibrium payoffs, although the bound will be weaker than Fudenberg and Levine’s. To see this suppose that the game is not of conflicting interests, but that there is a positive prior probability of a type who is committed to hold player two down to his minimax payoff. If the normal type of player one mimicks this type, then she can guarantee herself on average at least the payoff she would have got if she were publicly committed to this strategy. This doesn’t give her her most preferred payoff but it may still be more than her minimax payoff and thus reduce the set of Nash equilibrium payoffs as compared to the Folk theorem. I am grateful to Drew Fudenberg for this observation. A companion paper will generalize and elaborate this idea.

while the monopolist has to decide whether to acquiesce ( $A$ ) or to fight ( $F$ ). Assume that the payoffs of the unperturbed stage game are given as follows:

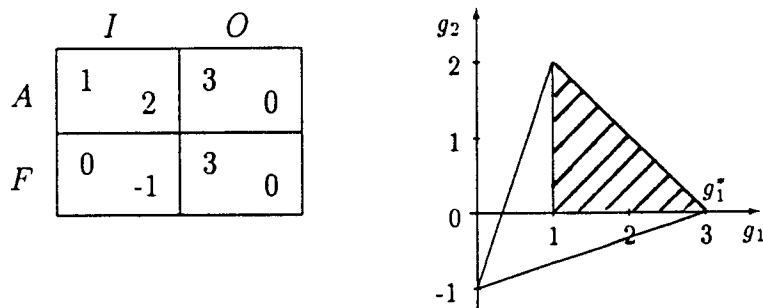


FIGURE 2 – The chain store game.

The monopolist would like to commit to fight in every period which would give her a commitment payoff of 3 and which would hold the entrant down to 0. Since 0 is also player two's minimax payoff the game is - according to our definition - of conflicting interests with respect to the monopolist. Kreps and Wilson (1982) have analyzed finite repetitions of this game with some incomplete information about the monopolist's type. For a particular perturbation of player one's payoff function they have shown that there are sequential equilibria in which the monopolist gets on average almost her commitment payoff if her discount factor is close enough to one and if there are enough repetitions. However, Fudenberg and Maskin (1986) demonstrated that any feasible payoff vector which gives each player more than his minimax payoff, i.e. any point in the shaded area of figure 2, can be sustained as an equilibrium outcome if the "right" perturbation is chosen. Thus, our Theorem 2 considerably strengthens the result of Kreps and Wilson (1982). It says that the only Nash equilibrium outcome of this game which is robust against any perturbation gives the monopolist at least her commitment payoff of 3 (note that she cannot get more), if she is sufficiently patient as compared to the entrant.<sup>8</sup>

<sup>8</sup>I am grateful to Eric van Damme for the following observation: Theorem 2 does not imply that the average payoff of player two is 0. Recall that player one is more patient than player two. So it may be that in the beginning of the game, say until period  $L$ , she gets less than 3 and player two gets more than 0, but after period  $L$  payoffs are always (3,0). For player one the first  $L$  periods do not count very much because she is very patient, so her average payoff is 3. However, player two puts more weight on the first  $L$  periods and less on everything thereafter, so her average discounted payoff may be considerably bigger than 0.

Furthermore it shows that this result carries over to the infinitely repeated game.

Now suppose that there is also incomplete information about the payoff function of the entrant. He would like to commit to enter in every period which would give him a commitment payoff of 2 while it would hold the monopolist down to 1, her minimax payoff. So the game is also of conflicting interests with respect to the entrant and our theorem applies. If there is two-sided uncertainty Proposition 2 says that it all depends on the relative patience of the two players and the prior probability distribution. If player one is sufficiently more patient than player two and if the probability of the normal type of player two is close to one, then player one will get her commitment payoff in any Nash equilibrium, and vice versa.

## 5.2 A Repeated Bargaining Game

Consider a buyer ( $b$ ) and a seller ( $s$ ) who bargain repeatedly in every period on the sale of a perishable good. The valuation of the buyer is 1 and the production costs of the seller are 0. Suppose there is a sealed bid double auction in every period: Both players simultaneously submit bids  $p_b$  and  $p_s$ ,  $p_i \in \{\frac{1}{n}, \frac{2}{n}, \dots, \frac{n}{n}\}$ , and there is trade at price  $p = \frac{p_b + p_s}{2}$  if and only if  $p_b \geq p_s$ . Consider the commitment strategy of the buyer. She would like to commit herself to offer  $p_b^* = \frac{1}{n}$  in every period. The unique best reply of the seller is  $p_s = \frac{1}{n}$ , which gives him  $g_s^* = \frac{1}{n}$ , his minimax payoff. Suppose the payoff function of the buyer is perturbed such that with some positive probability she will always offer  $p_b^*$ . Then Theorem 2 applies and the buyer will get almost her commitment payoff of  $\frac{n-1}{n}$  on average in any Nash equilibrium if her discount factor is close to one.

Note however, that this example is not as clear-cut as the chain store game. We have to assume that there is a minimal bid  $\frac{1}{n} > 0$ . If the buyer could offer  $p_b = 0$  she could hold the seller down to a minimax payoff of 0. But if he gets 0 the seller is indifferent between all possible prices, so he might choose  $p_s > 0$  and we end up with no trade. The point is that bargaining over a pie of fixed size is not quite a game of conflicting interests. Some cooperation is needed to ensure that trade takes place at all.

In Schmidt (1990) we considered a more complex extensive form game of repeated bargaining with one-sided asymmetric information, which confirms the above result that



the informed player can use the incomplete information about his type to credibly threaten to accept only offers which are very favourable to him. There, however, we took a different approach and it is interesting to compare the two models. In Schmidt (1990) we did not allow for “all possible” but only for “natural” perturbations of player one’s payoff function, i.e. we assumed that there may be incomplete information about the seller’s costs,  $c \in [0, 1]$ . The problem is that in this case there is no commitment type since none of the possible types of the seller has a dominant strategy in the repeated game. However, if the game has a finite horizon it can be shown that in any sequential equilibrium satisfying a weak Markov property the seller with the highest possible type will accept an offer if and only if it covers at least his costs. This type plays the role of the commitment type who is mimicked by all the other possible types of the seller. We show that the buyer will try to test the seller’s type at most a fixed finite number of times and that this will happen only in the end of the game. Surprisingly (from the point of view of Theorem 2) we can show that the seller will get his commitment payoff even if he is much less patient than the buyer, so the relative discount factors are not crucial. Furthermore the bargaining game we consider there is not of conflicting interests.<sup>9</sup> There are common interests as well, because players have to cooperate to some extent in order to ensure that trade takes place.

### 5.3 Games with Common and Conflicting Interests

“Pure” conflicting interests are a polar case and in most economic applications there are both - common and conflicting - interests present. Consider for example the repeated prisoner’s dilemma depicted in Figure 3. In a formal sense this game is of conflicting interests, but our theorem has no bite. Given that player two takes a best response against her commitment action player one would like most to commit herself to play  $D$ (effect) in every period. This holds player two down to his minimax payoff, but it only gives player one her minimax payoff as well. So, trivially she will get at least her commitment payoff in every Nash equilibrium. In this game the problem is not to commit to hold player two down to his minimax payoff, but to commit to cooperate.

---

<sup>9</sup>Note that not all possible perturbations are allowed for. This is why conflicting interests are not a necessary condition for the result in Schmidt (1990).

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

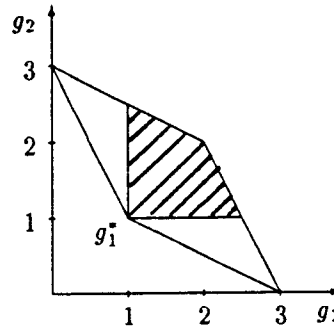


FIGURE 3 – The prisoner’s dilemma.

Another interesting example is a repeated Cournot game. Firm one would like to commit to choose the “Stackelberg-leader” quantity, which maximizes her stage game profit given that firm two chooses a best response against it. However, the “Stackelberg-follower” payoff of firm two is positive and thus greater than his minimax payoff since he can be held down to zero profits if player one gluts the market. So our result does not apply. Again the problem of player one is not to hold player two down as far as possible. Both players have the common interest to maximize joint profits, but interests are also conflicting in the sense that each of them would like to get more for himself at the expense of the other.

## 6 Extensions and Conclusions

To keep the argument as clear as possible we considered a very simple class of possible stage games with only two players, finite strategy sets, a countable set of possible types, and commitment types who always take the same pure action in every period. All of these assumptions can be relaxed without changing the qualitative results. Fudenberg and Levine (1989) provide a generalization to  $n$ -player games in which the strategy sets are compact metric spaces and in which there is a continuum of possible types of player one.<sup>10</sup> In Fudenberg and Levine (1992) they show that the argument can be extended to the case where the commitment types play mixed strategies and to games with moral

<sup>10</sup>If  $n \geq 3$ , the definition of a game of conflicting interests requires that  $a_1^*$  holds all other players  $i = 2, \dots, n$  down to their minimax payoffs simultaneously.

hazard, in which not the action of player one itself but only a noisy signal can be observed by player two. Since the technical problems involved in these generalizations are the same as in our model we refer to their work for any formal statements and proofs.

What happens if player one would like to commit herself to a more complex strategy which prescribes to take different actions in different periods and which may be conditional on player two's past play? It is easy to check that replacing  $a_1^*$  by  $a_1^*(h_t)$  doesn't change anything in the proofs of Lemmata 1 and 2 and Theorem 2. However, in a game with conflicting interests very little is gained by this generalization because  $a_1^*(h_t)$  has to hold player two down to her minimax payoff in any period and after any history  $h_t$ .

Fudenberg and Levine (1989) also demonstrated that the assumption that the stage game is simultaneous-move cannot be relaxed without an important qualification of their Theorem 1. The problem is that in an extensive form game player two may take an action after which player one has no opportunity to show that her strategy is the commitment strategy. Consider for example a repeated bargaining game in which in every period the buyer has to decide first whether to buy or not and then the seller has to choose whether to deliver high or low quality. If the buyer decides not to buy then he will not observe whether the seller would have produced high quality. This is why the seller might fail to get her commitment payoff in equilibrium. Note however that this problem does not arise in our context. The definition of a game with conflicting interests assumes that the commitment strategy of player one holds player two down to his minimax payoff. Therefore, if player two takes an action  $a_2$  in equilibrium after which player one's commitment strategy  $a_1^*$  is observationally equivalent to some other strategy  $a_1 \neq a_1^*$ , then player two cannot get more than his minimax payoff. So  $a_2$  must have been an element of  $B(a_1^*)$ . However player one's commitment payoff is defined as  $g_1^* = \max_{a_1 \in A_1} \min_{a_2 \in B(a_1^*)} g_1(a_1, a_2)$ . So if player two chooses  $a_2 \in B(a_1^*)$  player one cannot get less than  $g_1^*$ . Therefore, following Theorem 2 of Fudenberg and Levine (1989) it is straightforward that our result holds without qualification if  $g$  is any finite extensive form game.

To conclude, this paper has shown that "reputation effects" can explain commitment in a repeated game with two long-run players if and only if the game is of conflicting interests. If one of the players is very patient as compared to the other player, then

any Nash equilibrium outcome which is robust against perturbations of the information structure gives her on average almost her commitment payoff. This indicates that the message of the Folk theorem may be misleading. However, we still know very little about the evolution of commitment and cooperation in games in which both - common and conflicting - interests are present, which clearly is one of the most important issues of future research.

*Department of Economics, Massachusetts Institute of Technology, Cambridge, MA  
02139, U.S.A.*

## Appendix

### Proof of Lemma 2:

Consider any equilibrium  $(\sigma_1, \sigma_2)$  and fix a history  $h^t$  up to any period  $t$  along which player one has always played  $a_1^*$ , such that  $h^t$  has positive probability given  $(\sigma_1, \sigma_2)$ . Such a history exists because  $\mu^* > 0$ . Suppose that according to the (possibly mixed) equilibrium strategy  $\sigma_2^{t+1}$  player two chooses  $s_2^{t+1} \notin B(a_1^*)$  in period  $t+1$  with positive probability. Suppose further that the probability of player one not playing  $a_1^*$  (given that he always played  $a_1^*$  before) in each of the periods  $t+1, t+2, \dots, t+M$  is smaller than  $\epsilon$ . It will be shown that this can't be true in equilibrium because then player two would get less than his minimax payoff.

Note that  $\epsilon$  is independent of  $t$  and that  $M$  has been chosen in a way to guarantee that  $\epsilon > 0$ . Define  $\pi^\tau(a_1) = \text{Prob}(s_1^\tau = a_1 \mid h^{\tau-1})$  and let  $V_2^\tau(s_1^\tau, \sigma_2^\tau)$  be the continuation payoff of player two from period  $\tau$  onwards (and including period  $\tau$ ) given the strategy profile  $(s_1^\tau, \sigma_2^\tau)$  in period  $\tau$ . The expected payoff of player two from period  $t+1$  onwards is given by:

$$\begin{aligned}
 V_2^{t+1}(\sigma_1, \sigma_2) &= \sum_{a_1 \neq a_1^*} \pi^{t+1}(a_1) \cdot V_2^{t+1}(a_1, s_2^{t+1}) + \pi^{t+1}(a_1^*) \cdot \left\{ g_2(a_1^*, s_2^{t+1}) \right. \\
 &+ \delta_2 \cdot \sum_{a_1 \neq a_1^*} \pi^{t+2}(a_1) \cdot V_2^{t+2}(a_1, \sigma_2^{t+2}) + \delta_2 \cdot \pi^{t+2}(a_1^*) \cdot \left\{ g_2(a_1^*, \sigma_2^{t+2}) \right. \\
 (31) \quad &+ \dots + \\
 &+ \delta_2 \cdot \sum_{a_1 \neq a_1^*} \pi^{t+M}(a_1) \cdot V_2^{t+M}(a_1, \sigma_2^{t+M}) + \delta_2 \cdot \pi^{t+M}(a_1^*) \cdot \left\{ g_2(a_1^*, \sigma_2^{t+M}) \right. \\
 &\left. \left. + \delta_2 \cdot V_2^{t+M+1} \right\} \dots \right\}.
 \end{aligned}$$

It will be convenient to subtract  $\tilde{g}_2$  from both sides of the equation in every period. (Recall that  $\tilde{g}_2$  is the maximal payoff for player two if he takes an action which is not a best response against  $a_1^*$ .) Then we get:

$$\begin{aligned}
 V_2^{t+1}(\sigma_1, \sigma_2) - \frac{\tilde{g}_2}{1 - \delta_2} &= \sum_{a_1 \neq a_1^*} \pi^{t+1}(a_1) \cdot \left[ V_2^{t+1}(a_1, s_2^{t+1}) - \frac{\tilde{g}_2}{1 - \delta_2} \right] \\
 &+ \pi^{t+1}(a_1^*) \cdot \left\{ [g_2(a_1^*, s_2^{t+1}) - \tilde{g}_2] \right.
 \end{aligned}$$

$$\begin{aligned}
& + \delta_2 \cdot \sum_{a_1 \neq a_1^*} \pi^{t+2}(a_1) \cdot \left[ V_2^{t+2}(a_1, \sigma_2^{t+2}) - \frac{\tilde{g}_2}{1 - \delta_2} \right] \\
& + \delta_2 \cdot \pi^{t+2}(a_1^*) \cdot \left\{ [g_2(a_1^*, \sigma_2^{t+2}) - \tilde{g}_2] \right. \\
(32) \quad & + \dots + \\
& + \delta_2 \cdot \sum_{a_1 \neq a_1^*} \pi^{t+M}(a_1) \cdot \left[ V_2^{t+M}(a_1, \sigma_2^{t+M}) - \frac{\tilde{g}_2}{1 - \delta_2} \right] \\
& + \delta_2 \cdot \pi^{t+M}(a_1^*) \cdot \left\{ [g_2(a_1^*, \sigma_2^{t+M}) - \tilde{g}_2] \right. \\
& \left. + \delta_2 \cdot \left[ V_2^{t+M+1} - \frac{\tilde{g}_2}{1 - \delta_2} \right] \dots \right\} \left. \right\}.
\end{aligned}$$

By assumption the conditional probability that player one does not take her commitment action given that she always played  $a_1^*$  before is smaller than  $\epsilon$  in any period from  $t + 1, \dots, t + M$ , so

$$(33) \quad \sum_{a_1 \neq a_1^*} \pi^{t+i}(a_1) < \epsilon,$$

and, of course, we can use that  $\pi^{t+i}(a_1^*) \leq 1$ . Since  $\bar{g}_2$  is the maximal payoff player two can get at all, it has to be true that

$$(34) \quad V_2^{t+i}(a_1, \sigma_2^{t+i}) \leq \frac{\bar{g}_2}{1 - \delta_2} \text{ and } V_2^{t+M+1} \leq \frac{\bar{g}_2}{1 - \delta_2}.$$

Furthermore,  $s_2^{t+1}$  is supposed not to be a best response against  $a_1^*$ , so

$$(35) \quad g_2(a_1^*, s_2^{t+1}) \leq \tilde{g}_2.$$

Finally we can use that  $g_2(a_1^*, \sigma_2) \leq g_2^*$ . Substituting these expressions yields:

$$\begin{aligned}
& V_2^{t+1}(\sigma_1, \sigma_2) - \frac{\tilde{g}_2}{1 - \delta_2} < \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} + 1 \cdot \left\{ (\tilde{g}_2 - \tilde{g}_2) \right. \\
& \quad + \delta_2 \cdot \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} + \delta_2 \cdot 1 \cdot \left\{ (g_2^* - \tilde{g}_2) + \dots + \right. \\
& \quad \left. \left. + \delta_2 \cdot \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} + \delta_2 \cdot 1 \cdot \left\{ (g_2^* - \tilde{g}_2) + \delta_2 \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} \right\} \dots \right\} \right\}. \\
(36) \quad & = \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} + \delta_2 \cdot \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} + \delta_2 \cdot (g_2^* - \tilde{g}_2) \\
& \quad + \dots + \delta_2^{M-1} \cdot \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} + \delta_2^{M-1} \cdot (g_2^* - \tilde{g}_2) + \delta_2^M \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} \\
& = \epsilon \cdot (1 + \delta_2 + \dots + \delta_2^{M-1}) \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} + \delta_2^M \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2}
\end{aligned}$$

$$\begin{aligned}
& + (1 + \delta_2 + \dots + \delta_2^{M-1}) \cdot (g_2^* - \tilde{g}_2) - (g_2^* - \tilde{g}_2) \\
& < \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{(1 - \delta_2)^2} + \delta_2^M \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} - (g_2^* - \tilde{g}_2) + \frac{g_2^* - \tilde{g}_2}{1 - \delta_2}.
\end{aligned}$$

Recall from the statement of Lemma 2 that

$$(37) \quad \epsilon = \frac{(1 - \delta_2)^2 \cdot (g_2^* - \tilde{g}_2)}{\bar{g}_2 - \tilde{g}_2} - \delta_2^M \cdot (1 - \delta_2) > 0.$$

It is easy to check that  $\epsilon$  has been chosen such that

$$(38) \quad \epsilon \cdot \frac{\bar{g}_2 - \tilde{g}_2}{(1 - \delta_2)^2} + \delta_2^M \cdot \frac{\bar{g}_2 - \tilde{g}_2}{1 - \delta_2} = g_2^* - \tilde{g}_2.$$

Therefore we get:

$$(39) \quad V_2^{t+1}(\sigma_2) - \frac{\tilde{g}_2}{1 - \delta_2} < \frac{g_2^*}{1 - \delta_2} - \frac{\tilde{g}_2}{1 - \delta_2}.$$

However, since  $g_2^*$  is player two's minimax payoff this is a contradiction to the fact that we are in equilibrium. Q.E.D.

### Proof of Proposition 1:

The proof is similar to the construction of the counterexample in Section 3. Perturb the game  $g$  such that there are three types of player one, the normal type, the commitment type and an indifferent type, whose payoff is the same for any strategy profile, with probabilities  $(1 - \epsilon)$ ,  $\frac{\epsilon}{2}$ , and  $\frac{\epsilon}{2}$ , respectively. Let  $\underline{\delta}_2(\epsilon) = \frac{2}{2+\epsilon} < 1$  and suppose  $\delta_2 > \underline{\delta}_2(\epsilon)$ . Define

$$(40) \quad n = \frac{\ln \left[ 1 - \frac{2(1-\delta_2)}{\delta_2 \epsilon} \right]}{\ln \delta_2}$$

and let  $m = [n] + 2$ , where  $[n]$  is the integer part of  $n$ . Given the restriction on  $\delta_2$  it is straightforward to check that  $n$  is well defined and positive.

Since the commitment payoff of player one is strictly greater than her minimax payoff there exists an action  $\tilde{a}_2$  such that  $\tilde{g}_1 = g_1(a_1^*, \tilde{a}_2) < g_1^*$  and  $\tilde{g}_2 = g_2(a_1^*, \tilde{a}_2) < g_2^*$ . Suppose that  $\tilde{g}_1 > \min \max g_1$  and  $\tilde{g}_2 > \min \max g_2$ .<sup>11</sup> We will now construct an equilibrium such that the limit of the average equilibrium payoff of the normal type of player one for  $\delta_1 \rightarrow 1$  is bounded away from her commitment payoff by at least  $\eta$ , where

$$(41) \quad \eta = \frac{1}{m} \cdot [g_1^* - \tilde{g}_1] > 0.$$

<sup>11</sup>If for any of the players  $\tilde{g}_i \leq \min \max g_i$ , the construction of the "punishment equilibria" which are used below to deter any deviation from the equilibrium path are slightly more complex. In this case players have to alternate between the outcomes  $g^*$  and  $\tilde{g}$  such that both get on average at least their minimax payoffs.

Suppose  $1 > \delta_1 \geq \sqrt{\frac{\bar{g}_1 - \tilde{g}_1}{\bar{g}_1 - \tilde{g}_1 + g_1^* - \tilde{g}_1}}$ , where  $\bar{g}_1$  is the maximum payoff player one can get at all. Along the equilibrium path all types of player one play  $a_1^*$  in every period, while player two plays  $a_2^* \in B(a_1^*)$  in the first  $m - 1$  periods, then he plays  $\tilde{a}_2$  in period  $n$ , then starts again playing  $a_1^*$  for the next  $m - 1$  periods and so on. If player one ever deviates from this equilibrium path player two believes that he faces the normal type with probability 1. In this case we are essentially back in a game with complete information where the Folk theorem tells us that any individually rational, feasible payoff vector can be sustained as a subgame perfect equilibrium. So without writing down the strategies explicitly we can construct a continuation equilibrium, such that the continuation payoff is  $(\frac{1}{1-\delta_1}\tilde{g}_1, \frac{1}{1-\delta_2}\tilde{g}_2)$ . Clearly, the commitment and the indifferent type of player one have no incentive to deviate since  $a_1^*$  is at least weakly dominant for both of them. It is easy to check that - given  $m \geq 2$  and the restriction on  $\delta_1$  - the normal type of player one will not deviate either.

Now suppose player two ever deviates in any period  $t$ . In this case the normal and the commitment type are supposed to play  $a_1^*$  in period  $t + 1$ , while the indifferent type switches to another strategy  $\tilde{s}_1^{t+1} \neq a_1^*$ . If player two does not observe  $a_1^*$  being played in period  $t + 1$  he puts probability one on the indifferent type. Using the Folk theorem we can construct a continuation equilibrium in this subform which gives player two  $\frac{1}{1-\delta_2}\tilde{g}_2$  and which would give the normal type of player one  $\frac{1}{1-\delta_1}\tilde{g}_1$ . If, however, player two observes  $a_1^*$  being played in period  $t + 1$  he puts probability 0 on the indifferent type. In the continuation equilibrium of this subform  $(a_1^*, a_2^*)$  are always played along the equilibrium path. If there is any deviation by player one, player two believes that he faces the normal type with probability one and - using the Folk theorem again - the continuation payoff is  $(\frac{1}{1-\delta_1}\tilde{g}_1, \frac{1}{1-\delta_2}\tilde{g}_2)$ . Clearly, always to play  $a_2^*$  is a best response of player two against always  $a_1^*$  and always  $a_1^*$  is a best response for the commitment type against any strategy. It is easy to check that it is also a best response for the normal type of player one, given the "punishment" after any deviation.

We have already shown that the strategies of the players form an equilibrium after any deviation from the equilibrium path and that given the continuation equilibria player one has no incentive to deviate from this path. We still have to check that player two's strategy is a best response along the equilibrium path. The best point in time for a



deviation is when player two is supposed to play  $\tilde{a}_2$ . If it does not pay to deviate in this period, it never will. Suppose player two does not deviate. Then his payoff is given by:

$$\begin{aligned}
(42) \quad V_2(\tilde{a}_2) &= \tilde{g}_2 + \sum_{t=1}^{m-1} \delta_2^t g_2^* + \delta_2^m \tilde{g}_2 + \sum_{t=m+1}^{2m-1} \delta_2^t g_2^* + \dots \\
&= \tilde{g}_2 + \frac{\delta_2}{1-\delta_2} \cdot g_2^* - \frac{\delta_2^m}{1-\delta_2^m} \cdot (g_2^* - \tilde{g}_2).
\end{aligned}$$

However, if he deviates, the best he can do is to play  $a_2^*$  in period  $t$ . In this case his payoff is given by

$$(43) \quad V_2(a_2^*) = g_2^* + \delta_2 \cdot \left\{ \left(1 - \frac{\epsilon}{2}\right) \cdot \frac{1}{1-\delta_2} \cdot g_2^* + \frac{\epsilon}{2} \cdot \frac{1}{1-\delta_2} \cdot \tilde{g}_2 \right\}$$

It is now easy to check that  $\epsilon$  and  $\underline{\delta}(\epsilon)$  have been constructed such that  $V_2(\tilde{a}_2) > V_2(a_2^*)$ . Thus we have established that this is indeed an equilibrium path.

We now have to show that along this equilibrium path the average payoff of the normal type of player one is indeed smaller than  $g_1^* - \eta$  when  $\delta_1 \rightarrow 1$ . The equilibrium payoff of the normal type is given by:

$$\begin{aligned}
(44) \quad V_1 &= \sum_{t=1}^{m-1} \delta_1^{t-1} g_1^* + \delta_1^{m-1} \tilde{g}_1 + \sum_{t=m+1}^{2m-1} \delta_1^{t-1} g_1^* + \dots \\
&= \frac{1}{1-\delta_1} \cdot g_1^* - \frac{1}{\delta_1} \cdot \frac{\delta_1^m}{1-\delta_1^m} \cdot [g_1^* - \tilde{g}_1].
\end{aligned}$$

Therefore the difference between her commitment payoff and her average payoff in this equilibrium is

$$\begin{aligned}
(45) \quad g_1^* - (1-\delta_1) \cdot V_1 &= g_1^* - g_1^* + \frac{1-\delta_1}{\delta_1} \cdot \frac{\delta_1^m}{1-\delta_1^m} \cdot [g_1^* - \tilde{g}_1] \\
&= \frac{(1-\delta_1) \cdot \delta_1^{m-1}}{1-\delta_1^m} \cdot [g_1^* - \tilde{g}_1] \\
&= \frac{(1-\delta_1) \cdot \delta_1^{m-1}}{(1-\delta_1) \cdot \frac{1-\delta_1^m}{1-\delta_1}} \cdot [g_1^* - \tilde{g}_1] \\
&= \frac{\delta_1^{m-1}}{\sum_{t=0}^{m-1} \delta_1^t} \cdot [g_1^* - \tilde{g}_1] > \frac{\delta_1^{m-1}}{m} \cdot [g_1^* - \tilde{g}_1].
\end{aligned}$$

Taking the limit for  $\delta_1 \rightarrow 1$  we get

$$(46) \quad \lim_{\delta_1 \rightarrow 1} \frac{\delta_1^{m-1}}{m} \cdot [g_1^* - \tilde{g}_1] = \frac{1}{m} \cdot [g_1^* - \tilde{g}_1] = \eta.$$

Q.E.D.

## References

- AUMANN, R. AND S. SORIN (1989): "Cooperation and Bounded Recall", *Games and Economic Behaviour*, Vol. 1, 5-39.
- CRIPPS, M. AND J. THOMAS (1991): "Learning and Reputation in Repeated Games of Incomplete Information" mimeo, University of Warwick, June 1991.
- FUDENBERG, D. (1990): "Explaining Cooperation and Commitment in Repeated Games", forthcoming in *Advances in Economic Theory, Sixth World Congress* ed. by J.-J. Laffont, Cambridge University Press.
- FUDENBERG, D., D. KREPS, AND D. MASKIN (1990): "Repeated Games with Short-Run and Long-Run Players", *Review of Economic Studies*, Vol. 57, 555-573.
- FUDENBERG, D. AND D.K. LEVINE (1989): "Reputation and Equilibrium Selection in Games with a Patient Player", *Econometrica*, Vol. 57, 759 - 778.
- FUDENBERG, D. AND D.K. LEVINE (1992): "Maintaining a Reputation when Strategies are Imperfectly Observed", *Review of Economic Studies*, forthcoming.
- FUDENBERG, D. AND E. MASKIN (1986): "The Folk Theorem in Repeated Games with Discounting or with Incomplete Information", *Econometrica*, Vol. 54, 533-554.
- HARSANYI, J. (1967-68): "Games with Incomplete Information Played by Bayesian Players", *Management Science*, Vol. 4, 159-182, 320-334.
- HART, S. (1985): "Nonzero-Sum Two-Person Repeated Games with Incomplete Information" *Mathematics of Operations Research*, Vol. 10, 117-153.
- KREPS, D., P. MILGROM, J. ROBERTS AND R. WILSON (1982): "Rational Cooperation in the Finitely Repeated Prisoner's Dilemma", *Journal of Economic Theory*, Vol. 27, 245-252.
- KREPS, D. AND R. WILSON (1982): "Reputation and Imperfect Information", *Journal of Economic Theory*, Vol. 27, 253 - 279.
- MILGROM P. AND J. ROBERTS (1982): "Predation, Reputation and Entry Deterrence", *Journal of Economic Theory*, Vol. 27, 280 - 312.
- SCHMIDT, K. (1990): "Commitment through Incomplete Information in a Simple Repeated Bargaining Game", Discussion Paper No. A-303, Bonn University.
- SELTEN, R. (1978): "The Chain-Store Paradox" *Theory and Decision*, Vol. 9, 127-159.

## Footnotes

1. This paper is based on Chapter 3 of my PhD thesis which was completed within the European Doctoral Programme at Bonn University. I would like to thank David Canning, In-Koo Cho, Benny Moldovanu, Georg Nöldeke, Ariel Rubinstein, Avner Shaked, Joel Sobel, Monika Schnitzer, Eric van Damme and in particular Drew Fudenberg for many helpful comments and discussions. Financial support by Deutsche Forschungsgemeinschaft, SFB 303 at Bonn University, is gratefully acknowledged.
2. See Section 6 for the extension to extensive form stage games, continuous strategy spaces and more than two players.
3. Fudenberg and Levine (1989) refer to  $g_1^*$  as the “Stackelberg payoff”. However, it is now customary to use this expression only for  $\max_{a_1} \max_{\alpha_2 \in B(a_1)} g_1(a_1, \alpha_2)$ , that is for the maximum payoff player one could get if he could publicly commit himself to any action  $a_1$  and player two chooses the best response player one prefers *most*. See Fudenberg (1990). The analysis can be extended to the more general case where player one would like to commit himself to a mixed strategy or to a strategy dependent on history. See Fudenberg and Levine (1992) and the remarks in Section 6.
4. Whether there exists a *sequential* equilibrium in which player one gets substantially less than 10 if there are only the normal and the commitment type around is an open question. Note, however, that we want to characterize equilibrium outcomes which are robust to general perturbations of the informational structure of the game. From this perspective it makes little sense to restrict attention to two possible types only.
5. It is well known that a player can always guarantee herself at least her minimax payoff in any Nash-equilibrium.
6. Fudenberg and Maskin’s Folk theorem for games with incomplete information considers only finitely repeated games without discounting. However, the extension to discounting and an infinite horizon is straightforward.
7. Note that even if the game is not of conflicting interests we can still impose some

restriction on the set of Nash equilibrium payoffs, although the bound will be weaker than Fudenberg and Levine's. To see this suppose that the game is not of conflicting interests, but that there is a positive prior probability of a type who is committed to hold player two down to his minimax payoff. If the normal type of player one mimicks this type, then she can guarantee herself on average at least the payoff she would have got if she were publicly committed to this strategy. This doesn't give her her most preferred payoff but it may still be more than her minimax payoff and thus reduce the set of Nash equilibrium payoffs as compared to the Folk theorem. I am grateful to Drew Fudenberg for this observation. A companion paper will generalize and elaborate this idea.

8. I am grateful to Eric van Damme for the following observation: Theorem 2 does not imply that the average payoff of player two is 0. Recall that player one is more patient than player two. So it may be that in the beginning of the game, say until period  $L$ , she gets less than 3 and player two gets more than 0, but after period  $L$  payoffs are always (3,0). For player one the first  $L$  periods do not count very much because she is very patient, so her average payoff is 3. However, player two puts more weight on the first  $L$  periods and less on everything thereafter, so her average discounted payoff may be considerably bigger than 0.
9. Note that not all possible perturbations are allowed for. This is why conflicting interests are not a necessary condition for the result in Schmidt (1990).
10. If  $n \geq 3$ , the definition of a game of conflicting interests requires that  $a_1^*$  holds all other players  $i = 2, \dots, n$  down to their minimax payoffs simultaneously.
11. If for any of the players  $\tilde{g}_i \leq \min \max g_i$  the construction of the "punishment equilibria" which are used below to deter any deviation from the equilibrium path are slightly more complex. In this case players have to alternate between the outcomes  $g^*$  and  $\tilde{g}$  such that both get on average at least their minimax payoffs.

		<i>L</i>		<i>R</i>
<i>U</i>	10	10	0	0
<i>D</i>	0	0	1	1

“normal” type  
 $\mu^0 = 0.8$

		<i>L</i>		<i>R</i>
<i>U</i>	10	10	10	0
<i>D</i>	0	0	1	1

“commitment” type  
 $\mu^* = 0.1$

		<i>L</i>		<i>R</i>
<i>U</i>	1	10	1	0
<i>D</i>	1	0	1	1

“indifferent” type  
 $\mu^i = 0.1$

FIGURE 1 – A game with common interests.

	<i>I</i>	<i>O</i>
<i>A</i>	1    2	3    0
<i>F</i>	0    -1	3    0

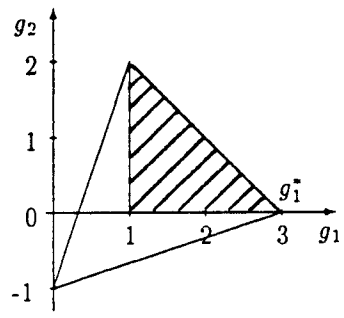


FIGURE 2 — The chain store game.

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

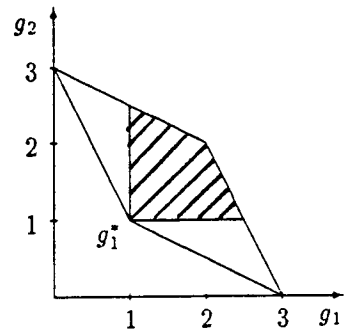


FIGURE 3 — The prisoner's dilemma.