# CONFLICT LEADS TO COOPERATION IN NASH BARGAINING

## By

## Kareen Rozen

**March 2008**
**Revised June 2009**

## COWLES FOUNDATION DISCUSSION PAPER NO. 1641

# Conflict Leads to Cooperation in Nash Bargaining*

Kareen Rozen[†]

Yale University

This Version: June 2009

## Abstract

We consider a multilateral Nash demand game where short-sighted players come to the bargaining table with requests for both coalition partners and the potentially generated resource. We prove that group learning leads with probability one to complete cooperation and a strictly self-enforcing allocation (i.e., in the interior of the core). Highlighting group dynamics, we demonstrate that behaviors which appear destructive can themselves lead to beneficial and strictly self-enforcing cooperation.

**Keywords:** Nash bargaining, learning, core, group conflict.
**JEL Codes:** C7.

# 1 Introduction

This paper develops a noncooperative model of multilateral bargaining in which group learning leads to convergence of allocations to the subset of the core which is *strictly* self-enforcing.

The literature on multilateral bargaining has a rich history, including Baron & Ferejohn (1989), Krishna & Serrano (1996), Chatterjee, Dutta, Ray & Sengupta (1993), Perry & Reny (1994), and Konishi & Ray (2003), with the latter three finding noncooperative bargaining foundations for the core (which consists of *weakly* self-enforcing allocations). In contrast to such papers, which typically examine the equilibria of dynamic bargaining games with forward-looking players, we are interested in the learning process resulting from repeated myopic play of a static bargaining game that extends the canonical two-player demand model of Nash (1950). We study learning in the bounded rationality or evolutionary sense of Fudenberg & Maskin (1992), Gale, Binmore & Samuelson (1995), and Mailath (1998), among others.

In our model, formalized in Section 2, $N$ players to come to the bargaining table with demands for both a potentially generated resource and coalition partners. The groups that form must be mutually compatible in terms of resource and partner requests, and their ability to produce the resource is governed by a convex and strictly superadditive characteristic function. We show in Section 3 that the set of strict Nash equilibrium outcomes of this static game correspond to the set of strictly self-enforcing (interior core) resource allocations and complete cooperation.

In Section 4 we study the learning process that results when this game is repeatedly played over time by myopic (or short-lived) players. By permitting players to include or exclude other players from their coalition, our model can capture flavorful and realistic group dynamics. In particular, group settings often display inefficient and destructive behaviors. Individuals can be excluded from groups or steal away other players' partners. Groups may take advantage of individuals who are desperate and alone, or have a scapegoat who absorbs the impact of a group failure. Individual greed may lead to internal strife, and one group's actions can instigate conflict within another. We construct a learning process by which the destructive behaviors of myopic individuals propel them towards strictly self-enforcing cooperation. Essentially, we show that the behaviors above are too destructive to sustain endless cycles of their use.

An interesting related paper by Agastya (1997) finds convergence to core (weakly self-enforcing) allocations through learning in a demand bargaining model where only resource requests are submitted. A main difference from this paper is that Agastya (1997) does not model coalition selection.[1] In our model, the dynamics of partner selection lay bare

---

[1]Because only resource demands are a strategic option for players, instead of modeling institutional details,

an interesting convergent process in which conflict itself leads to cooperation. Moreover, incorporating partner selection into the bargaining mechanism ensures that learning leads to a strictly self-enforcing outcome. We also show in Section 5.1 that our results extend when players have "pragmatic" preferences over partners: so long as excluding a given set of players would not affect her material payoff, the probability a player is willing to exclude them could depend on the properties of past play; for example, whether they excluded her earlier.

# 2 Nash bargaining with $N$ players

There is a group of $N \geq 3$ players, denoted by $I = \{1, 2, \ldots, N\}$. Letting $\mathcal{I}$ be the set of all possible coalitions, the resources a particular group may obtain is described by a convex and strictly superadditive characteristic function $v : \mathcal{I} \to \mathbb{R}$. Convexity means that there are increasing returns to scale: for all $S, T \subseteq I$, $v(S \cup T) - v(S) \geq v(T) - v(S \cap T)$. Strict superadditivity means that there are strictly positive synergies: if $S \cap T = \emptyset$, $v(S \cup T) > v(S) + v(T)$.

Players come to the bargaining table with two requests. First, as in the standard bilateral demand game, player $i$ requests some amount $d_i \in [v(i), v(I)]$ of the resource for herself. Second, player $i$ specifies a list of players $P_i \in \mathcal{I}$ with whom she is willing to form a coalition. For notational simplicity, we assume that player $i$'s list always includes herself. The list of all resource and partner requests submitted is given by $(d, P)$, where $d = (d_1, d_2, \ldots, d_N)$ and $P = (P_1, P_2, \ldots, P_N)$.

Not every combination of resource and partner requests is feasible. Letting $\Pi(I)$ denote the set of all coalition structures (i.e., partitions of $I$), a particular coalition structure $\pi \in \Pi(I)$ will be *feasible* if all of its coalitions are *mutually compatible* and *demand-feasible*. Mutual compatibility requires that for each group $S \in \pi$, no member $j \in S$ is excluded from the partner list of some other player in that group (i.e., there is no $i \in S$ such that $j \notin P_i$). Demand-feasibility is the simple condition that for each coalition $S \in \pi$ containing at least two players, the total amount of resource requested, $\sum_{i \in S} d_i$, does not exceed the

---

Agastya uses two characteristic-function based properties that determine whether or not a demand will be met. A limited form of coalition selection is permitted by Arnold & Schwalbe (2002), who allow players to switch only among existing coalitions (hence groups cannot split, and entirely new coalitions cannot be formed). They restrict the role of interaction by directly assuming that non-core allocations are unstable in Assumption 3, which says players in blocking coalitions may play randomly.

total amount of resource available, $v(S)$. *Strict demand feasibility* of the coalition $S$ means that $\sum_{i \in S} d_i < v(S)$.

Each player has a strictly increasing utility $u_i : \mathbb{R}_+ \to \mathbb{R}$ over the resource (for ease of exposition, we assume for now that the amount of resource received is the sole determinant of utility; in Section 5.1 we consider a class of preferences that permits some preferences for partners). A player $i$ who remains unpartnered in $\pi$ receives $v(i)$ regardless of her resource request. The utility of player $i$, under the requests $(d, P)$ and the coalition structure $\pi$, is given by $u_i(d_i)$ if $\pi$ specifies a nontrivial coalition for $i$, and $u_i(v(i))$ otherwise. A request $d_i$ is *individually rational* for player $i$ if $d_i \geq v(i)$, and *strictly individually rational* if the strict inequality holds. There is always a feasible and individually rational coalition structure: namely, the coalition structure where every player is unpartnered.

When more than one coalition structure is feasible, we assume that mutually compatible and demand-feasible groups form when possible. Formally, defining the norm $\rho : \Pi(I) \to \{1, 2, \ldots, N\}$ of a coalition structure to be the number of coalitions formed, we assume that the coalition structure that forms is chosen according to a fixed probability distribution with full support, $F \in \Delta\Pi(I)$, *conditional* on the set of feasible coalition structures with minimal $\rho$-norm. For example, if $N = 4$ and $P_i = \{1, 2, 3\}$ for $i = 1, 2$ and $P_j = \{1, 2, 3, 4\}$ for $j = 3, 4$, then the coalition structures of minimal $\rho$-norm are $\{(1, 2, 3), (4)\}$ and $\{(1, 2), (3, 4)\}$. Assuming these are both demand-feasible, which of these two coalition structure emerges will depend on the conditional distribution $F(\cdot \mid \{\{(1, 2, 3), (4)\}, \{(1, 2), (3, 4)\}\})$. Prior to knowing which $\pi$ forms, a player considers her $F$-expected utility over all feasible coalition structures of minimal $\rho$-norm.[2]

# 3 Enforceability through exclusion

The *core* of a cooperative game with characteristic function $v$, defined over the set of players $I$, is the set of self-enforcing allocations $\mathrm{Core}(v, I) = \left\{ d \mid \sum_{i \in I} d_i = v(I) \text{ and } \sum_{i \in S'} d_i \geq v(S') \ \forall \ S' \subset I \right\}$. We will be interested in the set of all strictly self-enforcing allocations (i.e., the *interior* of the core, obtained by using strict inequalities above), which we denote

---

[2]As noted in Hart & Kurz (1983), which considers coalition formation more generally, it is not evident how to predict which coalition structure forms if, for example, some member leaves a group. The minimal norm rule is meant to refine the prediction. The notions of coalitional compatibility suggested in Hart & Kurz (1983) are related to but differ from the definition here.

Core$^*(v, I)$. The set of such allocations is nonempty: convexity implies that the core is nonempty, and strict superadditivity implies it has a nonempty interior. In fact, each interior core allocation corresponds to a strict Nash equilibrium outcome of our demand game.

**Theorem 1** (Core equivalence). *$(d, P)$ is a strict Nash equilibrium outcome of the demand game if and only if $d \in Core^*(v, I)$ and $P_i = I$ for all $i$.*

That an interior core allocation and $P_i = I$ must be a strict Nash equilibrium outcome is clear: any deviation surely yields a player strictly less of the resource. To understand why being at an interior core allocation is necessary, it is helpful to make the following observations.

**Observation 1.** *In any strict Nash equilibrium, it must be that $P_i = I$ for all $i$, that the grand coalition is not strictly demand-feasible, and that $d_i > v(i)$ for all $i$.*

Indeed, switching from $P_i$ to any $P_i'$ with $P_i \subset P_i'$ does at least as well: either the resulting coalition structure has the same norm (in which case there is a weak increase in the probability that $i$ will be in a nontrivial coalition), or the norm decreases (in which case $i$ must have a partner, otherwise that coalition structure would have been feasible before). Moreover, even though the grand coalition must be mutually compatible, it cannot be strictly feasible because players would want to increase their resource requests.

Instead of concentrating on demand requests, our proof concentrates on when the players have disincentives to *exclude* others, building upon what we call the *exclusion principle*: *you should never exclude a player who can steal away members of your coalition and leave you alone.* Excluding such a player increases your probability of remaining unpartnered and receiving only $v(i)$, thereby lowering your expected utility. The following example illustrates.

**Example 1.** *Suppose that $I = \{1, 2, 3\}$ and that $v(i) = 0$ for all $i \in I$, $v(\{i, j\}) = 1$ for all $i \neq j$, and $v(I) = 3$. Then*

$$Core^*(v, I) = \{(d_1, d_2, d_3) \mid d_i > 0 \text{ for all } i, d_i + d_j > 1 \text{ for all } i \neq j, \text{ and } d_1 + d_2 + d_3 = 3\}.$$

*Consider the allocation $d = (\varepsilon, \varepsilon, 3 - 2\varepsilon) \notin Core^*(v, I)$ for any $\varepsilon \leq \frac{1}{2}$ (although it is in the core for $\varepsilon = \frac{1}{2}$) and note that this is not a strict Nash equilibrium allocation for any $\varepsilon \leq \frac{1}{2}$: player 1 can deviate by excluding player 3, since the coalition structure $(\{1\}, \{2, 3\})$ is infeasible by the resource constraint.*

The idea of the proof is as follows. If some player $i$ has a feasible subgroup, then players outside $i$'s subgroup must be able to steal away some of $i$'s partners to prevent $i$ from excluding them. That is, $i$ must have some chance of ending up alone in the resulting coalition structure in order to ensure that she strictly prefers not to restrict her set of acceptable partners. But then, the same disincentive to exclude must exist for that other feasible subgroup, and so and so forth. Continuing in this manner, more and more subgroups must be feasible until finally, by convexity, the grand coalition will itself be strictly feasible. However, if the grand coalition is strictly feasible, some player has an incentive to raise her resource request, which contradicts being at a strict Nash equilibrium.

*Proof of Theorem 1.* We now show that repeated use of the exclusion principle implies that $I$ cannot contain any demand-feasible subgroup. This would complete the proof, since being at a strict Nash equilibrium would then require that $\sum_{i \in \mathcal{I}} d_i = v(I)$. Suppose by contradiction that a demand-feasible subgroup does exist. If there is exactly one such subgroup, then any player $i$ inside it may exclude any player $j$ outside it (i.e., $j \notin P_i$) without affecting the feasible coalition structures containing $i$, and therefore without affecting $i$'s payoff — a contradiction to being at a strict Nash equilibrium. Therefore, there must be more than one feasible subgroup under $d$. Let $\hat{I}$ be the collection of players who have some feasible subgroup. We aim to show that $\hat{I}$ is strictly demand-feasible: for if $\hat{I} = I$, then the grand coalition is strictly feasible, and if $\hat{I} \subset I$, then the minimal norm rule ensures that any player within $\hat{I}$ may safely exclude any player outside $\hat{I}$.

Suppose that $\hat{I}$ is not feasible and that $S_1$, the largest feasible subgroup of $\hat{I}$, has size $s$. To prevent any player $i \in S_1$ from excluding any player $j \notin S_1$, it must be the case (by the exclusion principle) that $j$ must have a feasible subgroup $S_2$ containing some of $i$'s partners in $S_1$. For it to be possible that $i$ could remain alone in a a feasible coalition structure of minimal norm if she excludes $j$, it must be the case that $j$'s potential coalition $S_2$ also has size $s$. To see this, note that no subgroup strictly outside of $S_1$ can be feasible, else the union of the two would be feasible; and that for $i$ to remain alone, the norm cannot increase when $j$ steals $i$'s partners. That is, the minimal norm rule here means that the only way a player $i$ can steal the partners of another player $j$ is if player $i$ can command a coalition which is at least as large as $j$'s. The same exclusion principle holds for players in this next feasible subgroup $S_2$, and so on and so forth. Let $\{S_n\}_{1 \leq n \leq \hat{N}}$ denote the collection of all the feasible subgroups of size $s$. This collection must satisfy two properties:

(1) No player can be in every largest feasible subgroup (i.e., $\bigcap_{n \in \{1, \ldots, \hat{N}\}} S_n = \emptyset$).

(2) If $S_n \neq S_{n'}$, then $S_n \cap S_{n'}$ is a feasible subgroup.

Property (1) follows from the exclusion principle. Property (2) is a result of the following simple observation and the fact that $s$ is the size of the largest subgroup.

**Observation 2.** *Suppose that the resource request vector d has only strictly individually rational requests. If the two subgroups $S_n$ and $S_{n'}$ are demand-feasible and $S_n \cap S_{n'}$ is demand-infeasible or empty, then $S_n \cup S_{n'}$ is strictly demand-feasible.*

This observation follows directly from the definition of convexity. Notice that by property (2), $S_1 \cap S_2$ must be a feasible subgroup, else the subgroup $S_1 \cup S_2$ would be feasible and of size larger than $s$. Inductively, for every $k \leq \hat{N}$, $\cap_{j=1}^{k} S_j$ must be a feasible subgroup, else $S_k \cup \left( \cap_{j=1}^{k-1} S_j \right)$ would be feasible and of size larger than $s$. But then $\cap_{j=1}^{\hat{N}} S_j$ must be nonempty, contradicting property (1) and completing the proof. $\qquad\square$

# 4   Learning to cooperate from conflict

We now consider the $N$-player demand bargaining game played over time $t = 1, 2, \ldots$. As in much of the learning literature, the players can be interpreted as either successive generations of short-lived players or as a fixed set of myopic players.[3] The players respond only to the list of resource and partner requests $(d, P)$ submitted in the previous period. Typically (with probability $1 - \nu$ very close to one, and independently of other players), a player chooses a myopic best response to the previous period's demands. When there are multiple best responses, the player may choose any one of the strategies among which she is indifferent. With a small probability $\nu$, however, a player is inert: she does not update her request, leaving the previous period's demand in effect. Inertia may be interpreted in multiple ways; these include capturing exogenous constraints on the ability to actively bargain, difficulties in coordinating the timing of demands, learned behavior in the case of successive generations, or the manifestation of bounded rationality (e.g., attentional issues, computational costs, or simply slow updating of suboptimal strategies).

---

[3]The arguments also extend immediately to the case of multiple parallel populations that each population samples or more general matching technologies.

At each point in time, the previous period's requests $(d, P)$ serve as the state of the game. To ensure existence, player $i$'s resource requests are restricted to the discretized set $[v(i), v(I)]_K$ of $K$-place decimal fractions in $[v(i), v(I)]$.[4] The evolution of the game then defines a finite-state Markov chain over the state space of players' partner and resource requests. We are interested in how the group learns to play over time.

**Theorem 2** (Learning). *For sufficiently large $K$, the bargaining game converges with probability one to a state where $d \in \mathrm{Core}^*(v, I)$ and $P_i = I$ for all $i$.*

Clearly, any strict Nash equilibrium corresponds to an absorbing state of the dynamic process. Therefore, to prove this theorem we need only show that from any other state, the process can reach an interior core allocation with positive probability.

In particular, we show that there is positive probability that the following sequence of events will occur, in which eventual cooperation is the byproduct of familiar destructive behaviors. If players in a group cannot agree on an interior core allocation, then they may split into factions. Consequently, players may reach a situation where they are partitioned into mutually exclusive blocs, each of which agrees on an interior core allocation of their group. If any of these blocs consists of a lone player, then that player is desperate to receive any strictly individually rational amount and can offer to accept strictly less than her marginal contribution to some group. If that group takes advantage of her offer, an interior core allocation of the enlarged group can be created. With only nontrivial blocs remaining, each agreeing on an interior core allocation, one group $S$ instigates conflict over resources within another group $S'$ by inviting it to join and then rescinding the invitation - after the invitees have all responded greedily. With the abandoned group $S'$ unable to agree on a feasible allocation, one member is scapegoated and bears the burden of lowering her request. If this happens repeatedly, the scapegoat eventually leaves the group and can be picked up by $S$, creating an interior core allocation of the enlarged group. This process can then repeat itself until $S$ becomes the grand coalition. On the surface, these events take the appearance of a "divide and conquer" process, although the players involved are myopic.

We now develop this argument more formally.

*Proof of Theorem 2.* As a preliminary step in the proof, consider groups which are alienated from other players. Formally, suppose that the game is at a state $(d, P)$ where $d$ is not an

---

[4]We assume there is $K^*$ such that the values of $v$ are $K^*$-place decimal fractions and that $K \geq K^*$. It will be the case that best responses are always in $[v(i), v(I)]_K$.

interior core allocation and there exists a group of players $S \subseteq I$ such that every member of $S$ is excluded by every player outside $S$ ($P_i \cap S = \emptyset$ for all $i \notin S$), and vice-versa ($P_i \subseteq S$ for all $i \in S$). We first show there is a positive probability that either the players in $S$ come to agree on an allocation in the interior core of their group, or disintegrate into factions. To state this, we introduce the notation $d|_S$ and $v|_S$ for the restrictions of the allocation and characteristic function, respectively, to the group $S$. The resource allocation $d_S$ is *in the interior core of $S$* if $d|_S \in \text{Core}^*(v|_S, S)$.

**Lemma 1** (Factionization). *Suppose $d_S$ is not in the interior core of $S$ and $S$ is excluded by $I \setminus S$. Then there is positive probability that the game moves to a state $(d', P')$ where either the players in $S$ all agree to an allocation in the interior core of $S$ or a faction $T \subset S$ has broken away from $S$ (i.e., $P'_i = T$ for all $i \in T$).*

The proof of Lemma 1, which is in the appendix, builds on the exclusion technique developed earlier. If groups which cannot agree on an interior core allocation split into factions, then iterated application of Lemma 1 implies that from any nonabsorbing state, the game can transition within finite time to a state $(d^*, P^*)$ where the coalition structure is composed of mutually exclusive blocs, each in equilibrium with itself.

**Observation 3.** *It is possible to reach a coalition structure $\pi^*$ where every group is alienated from players outside it and agrees on an allocation in the interior core of their group (i.e., for all $S' \in \pi^*$, $P_i^* = S'$ for $i \in S'$, and if $S'$ is nonsingleton, then $d^*|_{S'} \in \text{Core}^*(v|_{S'}, S')$)*

If this coalition structure $\pi^*$ is the trivial one $\{(1), (2), \ldots, (N)\}$, then an interior core allocation is only a step away, for the players are indifferent among all requests. If $\pi^*$ is a nontrivial coalition structure then the situation is a bit trickier. However, using the following result we can assume that every bloc consists of at least two players. Indeed, suppose that some player $j$ is unpartnered, and therefore willing to accept any amount of resource larger than $v(j)$. We show that player $j$ can join an existing group $S$ - and create an interior core allocation for $S \cup \{j\}$ - by offering to accept strictly less than her marginal contribution.

**Lemma 2** (Enlarging a strictly self-enforcing agreement). *For large enough $K$, the game can reach a state $(\tilde{d}, \tilde{P})$ where $S$ and $j$ cooperate on an interior core allocation (i.e., $\tilde{P}_i = S \cup \{j\}$ for all $i \in S \cup \{j\}$ and $\tilde{d}|_{S \cup \{j\}} \in \text{Core}^*(v|_{S \cup \{j\}}, S \cup \{j\})$) and $(\tilde{d}, \tilde{P})$ is the same as $(d^*, P^*)$ for all other individuals.*

The states in Lemma 1 and Lemma 2 correspond to weak Nash equilibria. We now exhibit a path of play to a strict Nash equilibrium outcome (an interior core allocation) using destructive group behaviors. Indeed, supposing there exist two distinct blocs $S$ and $S'$ (otherwise the argument is complete), the actions of $S$ can lead to internal strife over resources within $S'$, permitting defectors from $S'$ to join $S$ à la Lemma 2.

To see how this may happen, suppose that the members of $S$ and $S'$ mutually invite each other; that is, simultaneously, every $i \in S \cup S'$ requests $(\tilde{d}_i, S \cup S')$. In the next period, if the players in $S$ are inert, while every player $j \in S'$ best responds with the request $(\tilde{d}_j + v(S \cup S') - v(S) - v(S'), S \cup S')$ (i.e., each attempts to grab all the remaining surplus), then the following result proves that continuing to invite members of $S'$ gives no additional expected utility to members of $S$. That is, the requests of the members of $S'$ have rendered those players useless to $S$.

**Lemma 3** (Disposability). *No member of $S'$ may feasibly join a coalition with any member of $S$.*

Suppose that the members of $S$ abandon the members of $S'$, as they are willing to do in light of Lemma 3. Specifically, suppose each $i \in S$ best responds with $(d_i^*, S)$ and that the members of $S'$ are inert. Since $S'$ had been at an interior core allocation, their members are unable to form any feasible coalitions with each other. In fact, if there is any player $k \in S'$ who is unable to obtain a payoff bigger than $v(k)$ by lowering her request, she may as well exit the coalition by setting $P_k = \{k\}$ and eventually join $S$ à la Lemma 2.

Otherwise, at least one of the members of $S'$ will need to lower her request. Let us consider what happens when this burden falls on one individual. Fix a *scapegoat* $j \in S'$ and suppose she is the only player in $S'$ to lower her request in the next period. Suppose that $j$ can obtain her best-response payoff by creating a coalition with just a subgroup of $S'$; then she may as well modify her resource request accordingly and set $P_j = S''$, where $S'' \subset S'$ is the smallest subgroup of $S'$ with which $j$ may obtain her best payoff. Note that the resulting allocation would be in Core*$(v|_{S''}, S'')$. This group $S''$ could safely break away in the next period, and $S' \setminus S''$ could then itself split or reach an interior core allocation as Lemma 1 prescribes.

If the scapegoat $j$ can only obtain a payoff larger than $v(j)$ by creating a coalition with the entire group $S'$, then the resulting allocation will be in Core*$(v|_{S'}, S')$. But now suppose the process repeats itself with the same scapegoat: $S$ and $S'$ mutually invite each other,

$S'$ responds greedily, $S$ abandons $S'$, and $j$ bears the burden of lowering her request. This need only be repeated a finite number of times before the scapegoat $j$'s best response is to break away - at which point she may join $S$ *à la* Lemma 2. Furthermore, $S$ can repeat this process against other groups until it grows to become the grand coalition and an interior core allocation is reached. □

# 5    Discussion

This paper demonstrates how inefficient group behaviors can propel groups toward strictly self-enforcing cooperative outcomes. In essence, we have shown that the destructive behaviors used here to achieve cooperation are too destructive to sustain endless cycles of their use. We discuss two extensions of the model below.

## 5.1    Introducing preferences for partners

For ease of exposition, we have assumed that a player's preferences depend on the resource only. However, these results easily generalize to a class of preferences that permits players to be pragmatically "behavioral." So long as excluding a player does not affect the expected amount of resource obtained, the probability of a player being willing to exclude any given set of players could, more realistically, be modeled to depend on the properties of past play; such as which players have excluded them earlier, whether their request was recently satisfied, and whether the player is sympathetic to someone who is unpartnered or instead attempts to "fit in" by excluding a player who has been excluded by others.

More formally, assume that the player may be described by a stochastic process over the set of all preference relations over $\mathcal{I}$, where every state has positive (but possibly negligible) probability of being reached from any state, and where the probability of transition may depend on the play of the game. Each player has a lexicographic preference, where she cares primarily about the amount of resource obtained, and secondarily about maximizing her preference over partners in her current state. For example, a vindictive player might, with high probability, strictly prefer to exclude players who have excluded her earlier – so long as doing so would not affect her materially. While vindictive behavior evidently enforces cooperation when players are forward-looking, the hope for cooperation might dim when players are both myopic and vindictive. On the contrary, both our results and the

convergent process exhibited carry through in this setting: the decision to include or exclude a player simply becomes a matter of strict preference.

## 5.2   Sharper long run prediction

To refine our prediction of interior core convergence further, we may introduce random perturbations, as in Kandori, Mailath & Rob (1993), to show that as these shocks become negligible, the outcomes persisting in the long run (i.e., *stochastically stable*, or in the support of the limiting stationary distribution) correspond to those allocations within the interior of the core that minimize the maximum individual wealth.[5]   This corresponds to a long run lexicographic social preference for strict enforceability (primarily) and wealth equity (secondarily).

---

[5]See the supplement posted on the author's website for the proof. Agastya (1999) has a result of the same spirit for the core rather than the interior core, using a different learning process. Both papers generalize the two-player result in Young (1993).

# Appendix

**Proof of Lemma 1.** Throughout, assume without loss that $d_i > v(i)$, $P_i = S$ for all $i \in S$ unless stated otherwise, $\sum_{i \in S} d_i \geq v(S)$, and $S$ and $I \setminus S$ mutually exclude each other.

Imagine first that $S$ contains no feasible subgroups. If $\sum_{j \in S} d_j = v(S)$, then condition (1) of the lemma is satisfied and the proof is complete. If instead $\sum_{j \in S} d_j > v(S)$ and no subgroups are feasible, then whenever only one individual $k$ best responds and the rest remain inert, one of three things may happen: (a) the resulting allocation could be in the interior core of the restricted game (again satisfying condition (1)), (b) the resulting allocation $(d, P)$ would no longer be strictly individually rational - then $P_k = \{k\}$ is a best response for $k$ and $P_i = S \setminus \{k\}$ becomes a best response for $i \in S \setminus \{k\}$ in the subsequent period (satisfying condition (2) of the lemma), or (c) the resulting allocation is strictly individually rational for players in $S$, $\sum_{j \in S} d_j \geq v(S)$, and some subgroup of $S$ is feasible. Consider the only nontrivial case, (c). Define the largest group size $s_d = \max_{\{T \subset S : \sum_{j \in T} d_j \leq v(T)\}} |T|$ and the collection $\mathcal{T}_d = \{ T \subset S \mid \sum_{i \in T} d_i \leq v(T) \text{ and } |T| = s_d \}$. There are two subcases.

**Case (i).** There is $T \in \mathcal{T}_d$ such that for all $i \in T$, $d_i$ is a best response to $(d, P)$. If a state satisfying condition (2) cannot be reached, no player $j \in T$ may be indifferent between $P_j = T$ and $P_j = S$: if $j$ best-responds with $(d_j, T)$ and all others play the same best response, then $(d_k, T)$ would be a best response for every $k \in T$ in the following period and condition (2) would be satisfied. So $(d_j, S)$ must be strictly preferred to $(d_j, T)$ for every $k \in T$; this implies that for each $j \in T$, there is a feasible group of size $s_d$ containing another member of $T$ but not $j$. No player in $T$ can be in every feasible group of size $s_d$ that contains a member of $T$; and the intersection of these groups of size $s_d$ must be feasible, else a bigger group is feasible. A contradiction can be found as in the proof of Theorem 1.

**Case (ii).** For all $T \in \mathcal{T}_d$, there is $i \in T$ such that $d_i$ is not a best response to $(d, P)$. For each $(d, P)$ we may partition the members of $S$ into the following three groups:

$$T_{(d,P)} = \{i \in S \mid d_i \text{ is a best response to } (d, P) \},$$
$$T_{(d,P)}^{+} = \{i \in S \setminus T_{(d,P)} \mid \text{ there is a best response } d_i^* \text{ to } (d, P) \text{ with } d_i^* > d_i\}, \text{ and}$$
$$T_{(d,P)}^{-} = \{i \in S \setminus (T_{(d,P)} \cup T_{(d,P)}^{+}) \mid \text{ there is a best response } d_i^* \text{ to } (d, P) \text{ with } d_i^* < d_i\}.$$

Beginning at state $(d, P)$, let all players in $T_{(d,P)} \cup T_{(d,P)}^{-}$ be inert and let all players in $T_{(d,P)}^{+}$ raise their requests. Call the resulting state $(d', P')$. If $T_{(d',P')}^{+} = \emptyset$, stop; otherwise this can be repeated a finite number of times until $T_{(\tilde{d},\tilde{P})}^{+} = \emptyset$ in the resulting state $(\tilde{d}, \tilde{P})$.

13

Suppose that a state satisfying condition (2) cannot be reached. The outcome of every player in $S$'s best response to $(\tilde{d}, \tilde{P})$ must be strictly individually rational, else some $k \in S$ could best respond by setting $P_k = \{k\}$ and a state satisfying condition (2) might result. Also, $T \in \mathcal{T}_{\tilde{d}} \Rightarrow T \nsubseteq T_{(\tilde{d}, \tilde{P})}$, otherwise one returns to Case (1). Therefore, $T_{(\tilde{d}, \tilde{P})}^{-} \neq \emptyset$.

The first task is to show that under $(\tilde{d}, \tilde{P})$ and the assumption that condition (2) cannot be satisfied, $S$ cannot have any feasible subgroups. Suppose that there is at least one feasible subgroup of $S$, and again denote by $s_{\tilde{d}}$ the size of the largest such subgroup. $T_{(\tilde{d}, \tilde{P})}^{+} = \emptyset$ and $T \in \mathcal{T}_{\tilde{d}} \Rightarrow T \nsubseteq T_{(\tilde{d}, \tilde{P})}$, so some $i \in T_{(\tilde{d}, \tilde{P})}^{-}$ must be both included and excluded from feasible subgroups of $S$ of size $s_{\tilde{d}}$. If she were never excluded, lowering her request would not be a best response. Note once more that no player can be in every feasible subgroup of size $s_{\tilde{d}}$ (because condition (2) cannot be satisfied), and that the intersection of any two such subgroups must be a feasible subgroup (because no subgroup of size larger than $s_{\tilde{d}}$ is feasible and $S$ is not strictly feasible). The same argument as in Case (1) leads to the desired contradiction.

Hence, $S$ lacks feasible subgroups under $(\tilde{d}, \tilde{P})$. Choose some $k \in T_{(\tilde{d}, \tilde{P})}^{-}$ to best respond and let all others be inert. The best response of $k$ has $d_k^* = \max_{T \subseteq S, k \in T} v(T) - \sum_{j \in T \setminus \{k\}} \tilde{d}_j$. If $T^* \in \arg\max_{T \subseteq S, k \in T} v(T) - \sum_{j \in T \setminus \{k\}} \tilde{d}_j$ for some $T^* \neq S$, then $(d_k^*, T^*)$ is optimal for $k$. Next period, $(\tilde{d}_j, T^*)$ will be a best response for each $j \in T^* \setminus \{k\}$, a contradiction to the assumption that condition (2) cannot be satisfied. Therefore, $S$ forms and no subgroups of $S$ will be feasible, i.e. a state satisfying (1) will be reached. $\qquad\square$

**Proof of Lemma 2.** Any request is a best response for $j$; and those in $S$ are indifferent about inviting players who exclude them. Fix $m \in \mathbf{Z}_+$. Suppose that in the same period, player $j$ requests $(v(S \cup \{j\}) - v(S) - m \cdot 10^{-K}, S \cup \{j\})$ and each $i \in S$ requests $(d_i, S \cup \{j\})$. Next period, some $k \in S$ requests $(d_k + m \cdot 10^{-K}, S \cup \{j\})$ and players in $(S \cup \{j\}) \setminus \{k\}$ don't move. It remains to verify that the resulting allocation $d'$ has $d'|_{S \cup \{j\}} \in \text{Core}^*(v|_{S \cup \{j\}})$.

Clearly $\sum_{i \in S} d_i = v(S)$. Define $\varepsilon^* = \min_{S \cap T = \emptyset} v(S \cup T) - v(S) - v(T)$, which is positive by strict superadditivity, and assume $K$ is large enough that $m \cdot 10^{-K} < \varepsilon^*$. The assumption on $K$ guarantees that $d'_j > v(j)$ is satisfied. For any $S' \subset S$, we must show $S' \cup \{j\}$ is infeasible. If $k \in S'$, this is trivial by convexity and the fact that $d|_S \in \text{Core}^*(v|_S)$:

$$\sum_{i \in S' \cup \{j\}} d'_i = \sum_{i \in S'} d_i + v(S \cup \{j\}) - v(S) > v(S') + v(S \cup \{j\}) - v(S) \geq v(S') + v(S' \cup \{j\}) - v(S')$$

If $k \notin S'$, then the infeasibility requirement is satisfied when $\sum_{i \in S'} d_i > v(S') + m \cdot 10^{-K}$; for then convexity and $d'_i = d_i$ for $i \in S'$ ensure that $v(S' \cup \{j\}) < d'_j + \sum_{i \in S'} d'_i$ because

$$m \cdot 10^{-K} < v(S' \cup \{j\}) - v(S') + \sum_{i \in S'} d_i - v(S' \cup \{j\}) \leq v(S \cup \{j\}) - v(S) + \sum_{i \in S'} d_i - v(S' \cup \{j\})$$

A technical issue arises only when $\hat{S} = \left\{ S' \subset S, S' \neq \emptyset \mid \sum_{i \in S'} d_i \leq v(S') + m \cdot 10^{-K} \right\}$ is nonempty and such that $\cap_{S' \in \hat{S}} S' = \emptyset$. If $\cap_{S' \in \hat{S}} S' \neq \emptyset$, simply let the best-responding player $k$ be in $\cap_{S' \in \hat{S}} S'$. We will now show that $\exists\, K^{**} \in \mathbf{Z}_+$ such that $\cap_{S' \in \hat{S}} S' = \emptyset$ is impossible whenever $K \geq \max \left\{ K^*, K^{**} \right\}$. Let $K^{**} = [\log \frac{|\hat{S}|(m+1)+m}{\varepsilon^*}]+1$ and suppose that $\cap_{S' \in \hat{S}} S' = \emptyset$. Convexity necessitates that $\sum_{i \in S' \cap S''} d_i \leq v(S' \cap S'') + (2m+1) \cdot 10^{-K}$, otherwise $S' \cup S''$ is strictly feasible, a contradiction to $d|_S \in \mathrm{Core}^*(v|_S)$. Consider some $S' \in \hat{S}$ and take $S'' \subset S$ such that $S' \not\subseteq S''$ and $S'' \not\subseteq S'$. If $\sum_{i \in S''} d_i = v(S'') + r \cdot 10^{-K}$ for some $r \leq |\hat{S}|(m+1) + m$, then by convexity it must be that $\sum_{i \in S' \cap S''} d_i \leq v(S' \cap S'') + (r+m+1) \cdot 10^{-K}$ to avoid the contradiction that $S' \cup S''$ is strictly feasible. Consider two distinct $S_1, S_2 \in \hat{S}$ and let $T_1 = S_1 \cap S_2$. $T_1 \neq \emptyset$, else $S_1 \cup S_2$ is strictly feasible. If $T_1 \neq S_1, S_2$ then $\sum_{i \in T_1} d_i \leq v(T_1) + (2m+1) \cdot 10^{-K}$; otherwise, if $T_1 = S_1$ then $\sum_{i \in T_1} d_i = \sum_{i \in S_1} d_i$, and similarly for the case $T_1 = S_2$. In either case, $\sum_{i \in T_1} d_i \leq v(T_1) + (2m+1) \cdot 10^{-K}$ is the upper bound of interest. Inductively define $T_n = T_{n-1} \cap S_{n+1}$ for $2 \leq n \leq R = |\hat{S}| - 1$. If $T_n = \emptyset$, one obtains a contradiction. We are concerned with the case $T_n \neq T_{n-1}, S_{n+1}$ to get the upper bound on $\sum_{i \in T_n} d_i$, which by convexity is $\sum_{i \in T_n} d_i \leq v(T_n) + [(n+2)m+n+1] \cdot 10^{-K}$. The final intersection $T_R = T_{R-1} \cap S_{R+1}$ must be empty and a contradiction arises. $\qquad\square$

**Proof of Lemma 3.** Define $d'$ by $d'_i = \tilde{d}_i + [v(S \cup S') - v(S) - v(S')] \cdot 1_{i \in S'}$, where $1_X$ is the usual indicator function. First, we prove the following intermediate result using convexity: *take nonempty $A, A', B \subset I$ with $A \cap B = \emptyset$ and $A' \subset A$; and let $d$ be such that $d|_A \in \mathrm{Core}^*(v|_A)$. If $A' \cup B$ is a feasible coalition under $d$, then $A \cup B$ is strictly feasible under $d$.* To see this, note that by convexity, $v(A \cup B) - v(A) \geq v(A' \cup B) - v(A')$. By assumption, both $\sum_{i \in A'} d_i > v(A')$ and $\sum_{i \in A' \cup B} d_i \leq v(A' \cup B)$. Hence $v(A \cup B) - v(A) > \sum_{i \in B} d_i$. Noting that $\sum_{i \in A} d_i = v(A)$ completes the proof of the claim.

15

We claim that for any $\emptyset \neq S'' \subseteq S'$, $\sum_{i \in S''} d_i' > v(S \cup S'') - v(S)$. To see this, note that by convexity, strict superadditivity, and $\tilde{d}|_{S'} \in \text{Core}^*(v|_{S'})$,

$$\sum_{i \in S''} \tilde{d}_i + |S''|[v(S \cup S') - v(S) - v(S')] - v(S \cup S'') + v(S)$$

$$\geq \sum_{i \in S''} \tilde{d}_i + |S''|[v(S \cup S'') - v(S) - v(S'')] - v(S \cup S'') + v(S)$$

$$\geq (|S''| - 1)[v(S \cup S'') - v(S) - v(S'')].$$

If $S''$ is non-singleton the last term is strictly positive, and if $S''$ is singleton the intermediate term is strictly positive by strict individual rationality of the request. The lemma then follows from the contrapositive of the intermediate claim. $\qquad\square$

# References

Agastya, M. 1997. "Adaptive Play in Multiplayer Bargaining Situations." *Review of Economic Studies* 64:411–426.

Agastya, M. 1999. "Perturbed Adaptive Dynamics in Coalition Formation Games." *Journal of Economic Theory* 89:207–233.

Arnold, T. & U. Schwalbe. 2002. "Dynamic Coalition Formation and the Core." *Journal of Economic Behavior and Organization* 49:363–380.

Baron, D.P. & J. Ferejohn. 1989. "Bargaining in Legislatures." *American Political Science Review* 83:1181–1206.

Chatterjee, K., B. Dutta, D. Ray & K. Sengupta. 1993. "A Noncooperative Theory of Coalitional Bargaining." *Review of Economic Studies* 60:463–477.

Fudenberg, D. & E. Maskin. 1992. "Evolution and Cooperation in Noisy Repeated Games." *The American Economic Review, Papers and Proceedings* pp. 274–279.

Gale, J., K. Binmore & L. Samuelson. 1995. "Learning to Be Imperfect: The Ultimatum Game." *Games and Economic Behavior* 8:56–90.

Hart, S. & M. Kurz. 1983. "Endogenous Formation of Coalitions." *Econometrica* 51:1047–1064.

Kandori, M., G. Mailath & R. Rob. 1993. "Learning, Mutation and Long Run Equilibria in Games." *Econometrica* 61:29–56.

Konishi, H. & D. Ray. 2003. "Coalition Formation as a Dynamic Process." *Journal of Economic Theory* 110:1–41.

Krishna, V. & R. Serrano. 1996. "Multilateral Bargaining." *Review of Economic Studies* 63:61–80.

Mailath, G. 1998. "Do People Play Nash Equilibrium? Lessons from Evolutionary Game Theory." *Journal of Economic Literature* 36:1347–1374.

Nash, J. 1950. "The Bargaining Problem." *Econometrica* 18:155–162.

Perry, M. & P. Reny. 1994. "A Noncooperative View of Coalition Formation and the Core." *Econometrica* 62:795–817.

Young, H.P. 1993. "An Evolutionary Model of Bargaining." *Journal of Economic Theory* 59:145–168.