

---

**Evolutionary Bargaining with Cooperative Investments**

**Herbert Dawid and W. Bentley MacLeod**

**USC Center for Law, Economics & Organization  
Research Paper No. C02-19**



**CENTER FOR LAW, ECONOMICS  
AND ORGANIZATION  
RESEARCH PAPER SERIES**

**Sponsored by the John M. Olin Foundation**

University of Southern California Law School  
Los Angeles, CA 90089-0071

*This paper can be downloaded without charge from the Social Science Research Network  
electronic library at [http://papers.ssrn.com/abstract\\_id=xxxxxx](http://papers.ssrn.com/abstract_id=xxxxxx)*

---

# Evolutionary Bargaining with Cooperative Investments\*

Herbert Dawid

W. Bentley MacLeod

Department of Economics

University of Southern California

Los Angeles CA 90089-0253

January 31, 2002

## Abstract

This paper explores the set of stochastically stable equilibria in a model in which individuals first decide to make a high or low investment, and then are matched to play a Nash demand game. If an agreement is not reached, then they are re-matched in the next period, and obtain a payoff discounted by  $\delta$ . We identify a condition under which stochastically stable bargaining conventions exist and find, that the stochastically stable division rule is independent of the long run investment strategy. In these conventions the potential to trade in subsequent periods always has an effect on the bargain, and the market acts more like a threat point, than an outside option. If investments are substitutes stochastically stable bargaining conventions imply larger investment incentives than the Nash bargaining solution whereas the opposite is true if investments are complements. Finally, if it is not efficient for trade to occur as a result of the outside option, and investments are complements, then no bargaining convention can develop, and investment levels are typically inefficient.

---

\*We would like to thank Jack Robles for helpful comments and the National Science Foundation for financial support under grant SES 0095606.

# 1 Introduction

Efficient exchange often entails the use of relationship specific investments, but in the absence of binding contracts, the *ex post* negotiation of the terms of trade can result in the sharing of the gains from investment between the two parties, leading to the well know problem of holdup (Grout (1984) and Grossman and Hart (1986)). In the case of one-sided relationship specific investment, followed by the play of a Nash demand game to determine the terms of trade Tröger (2000) and Ellingsen and Robles (2000) find that under the appropriate conditions stochastically stable equilibria entail efficient investment. These results are quite surprising because they illustrate a situation in which learning, as modeled by the criteria of stochastic stability, leads to behavior that is sensitive to sunk costs. The purpose of this paper is to extend this work in two directions. First we allow for investment by both parties, and, second, individuals that do not trade can at a cost rematch and attempt trade with a new party in the subsequent period.

In our model whenever investment by both parties is efficient, there are a large number of efficient subgame perfect equilibrium for the stage game. Using the methodology of Young (1993) and Kandori, Mailath, and Rob (1993) we explore the set of stochastically stable equilibria in a model in which individuals first decide to make a high or low investment, and then are matched to play a Nash demand game. If an agreement is not reached they are re-matched in the next period, and then obtain a payoff discounted by  $\delta$ . We find, in contrast to Tröger (2000) and Ellingsen and Robles (2000), that in the case of two-sided investment the stochastically stable division rule in general does not provide efficient investment incentives, and hence holdup is still a problem. The potential to trade in subsequent periods always has an effect on the bargain for all  $\delta > 0$ , and therefore the market acts more like a threat point, than an outside option in the sense of Binmore, Rubinstein, and Wolinsky (1986). It turns out that for the allocation of surplus in the stochastically stable convention the value of the outside options in environments of low investment is crucial even if the induced long run outcome is full investment. This implies that when investments are substitutes then the set of parameter values yielding high investment is larger than in the standard holdup problem where the allocation follows the Nash bargaining solution. Conversely, when investments are complements, the criterion of stochastic stability makes the holdup problem worse. If it is not efficient for trade to occur as a result of the outside option, and investments are complements, then no bargaining convention can develop, and investment levels are typically inefficient.

The agenda of the paper is as follows. The next section introduces the basic model, followed by an illustration of the potential for efficient bargaining norms in this model. Sections 4 and 5 introduce the formal stochastic learning model, and present a preliminary analysis of the stochastically stable sets. Section 6 considers the case of substitutes, where the marginal return from the first investment is greater than the second investment, while section 7 presents our results for complementary investments. The paper concludes with a discussion of the results and their relationship to the literature.

## 2 The Model

We are interested in the kind of bargaining and investment conventions which are developed endogenously in a population of adaptive agents. To examine this, we use an evolutionary bargaining model similar to Young (1993) as extended to incorporate investment by Tröger (2000) and Ellingsen and Robles (2000). It is assumed that agents use a random sample of the population of players to build beliefs about the investment and bargaining behavior of the other individuals. With a large probability they then choose the optimal strategy given their beliefs.

Consider a single population of identical agents who are repeatedly matched randomly in pairs to engage in joint production (or in a joint project). Every agent can make an investment, either high (H) or low (L), before entering the population that influences his type, and accordingly the joint surplus of the project. This investment can be thought of as human capital, such as the acquisition of special skills needed for a project, though the framework is sufficiently general that any type of project specific investment might be considered. Before partners start joint production they bargain over the allocation of the joint surplus. If the bargaining does not lead to an agreement they split without carrying out the project and look for new partners. The effect of an investment stays intact as long as the agent has not carried out the project, it is however assumed that the investment is project specific and creates no additional revenue after the project has been carried out. The degree of project specificity is parameterized by a discount factor  $\delta \in [0, \bar{\delta}]$ , such that the value of trade  $t$  periods after the initial investment is  $\delta^t U$ , where  $U$  is the agents share of the gains from trade. When  $\delta = 0$  the model corresponds to purely relationship specific investment.

The sequence of decisions facing an individual are:

1. The agent,  $i$ , decides about her investment level  $I^i \in \{h, l\}$ , where the cost of investment is

$$c(I) = \begin{cases} c, & \text{if } I = h, \\ 0, & \text{if } I = l. \end{cases}$$

After the investment has been made the type  $T^i \in \{H, L\}$  of the agent is determined. It is assumed that the probability of being a high type after having invested  $I$  is  $p_I$ , where  $p_h > p_l$ .

2. The agent is randomly matched with some partner and both observe each other's type. The types determine the size of the surplus,  $S_{T_i T_j}$ , which satisfies  $S_{HH} \geq S_{LH} = S_{HL} \geq S_{LL} > 0$ .
3. Individual  $i$  makes a demand conditional upon her type and that of her partner  $j$ , denoted by  $x_{T_i T_j} \in X_{T_i T_j}(k) = \{0, \alpha_{T_i T_j}, 2\alpha_{T_i T_j}, \dots, k\alpha_{T_i T_j}\}$ ,  $\alpha_{T_i T_j} = S_{T_i T_j}/k$ ,  $k$  is some large even number.
4. The payoff to individual  $i$  in this period is given by the rules of the Nash demand game:

$$U^i = \begin{cases} x_{T_i T_j}^i, & \text{if } x_{T_i T_j}^i + x_{T_j T_i}^j \leq S_{I_i I_j} \\ 0, & \text{if } x_{T_i T_j}^i + x_{T_j T_i}^j > S_{I_i I_j} \end{cases} - c(I^i)$$

and similarly for player  $j$ . Agents are assumed to be risk neutral.

5. If agent  $i$  has traded in this period she leaves the population and is replaced by another individual. If there was no trade the individual stays in the population and goes again through steps 2 - 5 in the following period where future payoffs are discounted by a factor  $\delta$  per period.

Throughout the analysis  $S_{HH}$  and  $S_{LL}$  are assumed fixed, while the degree of complementarity in investment,  $S_{LH}$ , and the cost of investment,  $c$ , are parameters that determine the nature of the investment problem. Furthermore, we assume that the probability that the type differs from the investment level is symmetric and small, namely:  $1 - p_h = p_l = \lambda$  for some small positive  $\lambda$ . This latter assumption plays an important role in the analysis because it ensures that even if all individuals carry out high investment, there is a strictly positive probability of having low types in the population. Hence each period there is the potential for trade between  $H$  and  $L$  types. As we shall see, the existence of such trades is a necessary condition for the evolution of a bargaining convention.

From the analysis of Young (1993), it is known that the equal split is stochastically stable when all individuals are the same. Therefore, to simplify the analysis it is assumed that when two high types meet or two low types meet they split the gains from trade equally if they trade, i.e.  $x_{HH}^i = \frac{S_{HH}}{2}, x_{LL}^i = \frac{S_{LL}}{2} \forall i$ . For most of the current analysis it shall be assumed that the discount factor  $\delta$  is sufficiently small that it is always efficient to trade, regardless of the type of your partner, rather than wait. Hence the option to wait will act as a constraint on the current trade, an assumption that is discussed in more detail in the next section.

These assumptions greatly simplify the strategy space. When a player first enters the game she chooses  $I \in \{h, l\}$ , after which point she learns her type  $T \in \{H, L\}$ . Given her type, each period she needs to formulate only her demand when faced with a partner of a different type, since she adopts the equal split rule when faced with a partner of the same type. Formally, a strategy of the stage game is given by  $(I, x_{HL}, x_{LH}) \in \{h, l\} \times X(k)^2$ , where  $X(k) = X_{LH}(k) = X_{HL}(k)$ , but in every period other than the period she enters an agent only has to determine one action, namely  $x_{HL}$  if she is of type  $H$  or  $x_{LH}$  if she is of type  $L$ . In what follows we will refer to the pair  $(x_{HL}, x_{LH})$  as the bargaining strategy of an agent.

### 3 Equilibrium Analysis

Our goal is to understand the structure of the stochastically stable equilibria as a function of the cost of investment,  $c$ , the degree of investment complementarity,  $S_{LH}$ , and the degree of investment specificity, modeled by  $\delta$ . The purpose of this section is to characterize the uniform subgame perfect equilibria in stationary strategies of the population game<sup>1</sup> that result in high investment. It will turn out that if stochastically stable equilibria exist they are indeed in this class of equilibria.

Note that in the Nash demand game any strategy profile  $(x_{HL}, x_{LH})$  such that  $x_{LH} + x_{HL} = S_{LH}$  is a Nash equilibrium. By a *bargaining convention* we mean a situation where all individuals have identical bargaining strategies of the form  $(S_{LH} - \hat{x}_{LH}, \hat{x}_{LH})$  for some  $\hat{x}_{LH} \in [0, S_{LH}]$ .

Since the focus of this paper lies on the bargaining behavior in matches of different types we will make assumptions that guarantee that equal split trades always occur between equal types. Given the results of Young (1993) we only have to be concerned about the question whether equal types want to trade at all or rather wait for a different type. The maximal payoff a low type can get in the next period is  $S_{LH}$  and

<sup>1</sup>This means that we consider scenarios where all individuals use identical strategies of the stage game every period and these strategies are constant over time.

therefore  $\frac{S_{LL}}{2} > \delta S_{LH}$  is sufficient to guarantee trade between low types. For high types we must have  $\frac{S_{HH}}{2} > \delta S_{LH}$  which clearly is a weaker condition. Hence we will assume throughout the paper that

$$(1) \quad \delta < \frac{S_{LL}}{2S_{LH}}.$$

Considering High-Low pairings we realize that for relatively high discount factors and strong complementarity between investments, even if a bargaining convention exists, one of the two partners would rather wait for a partner of identical type than to trade according to the bargaining convention. Given that in a High-Low pairing the high type expects a low bid of  $\hat{x}_{LH}$ , the low type expects a high bid of  $S_{LH} - \hat{x}_{LH}$  and given that both partners believe that they will meet an identical type in the following period, they will be willing to trade if

$$\begin{aligned} S_{LH} - \hat{x}_{LH} &> \delta S_{HH}/2, \\ \hat{x}_{LH} &> \delta S_{LL}/2. \end{aligned}$$

The first condition ensures that the high type prefers trading with a low type, rather than waiting one period and trading with a high type. The second condition is the corresponding requirement for the low type. Adding these inequalities together implies the following necessary condition for trade to occur for HL matches:

$$(2) \quad \frac{2S_{LH}}{S_{LL} + S_{HH}} > \delta.$$

Put differently, (2) implies that there exists a bargaining convention  $x_{LH}$  such that individuals always trade in High-Low matchings no matter what their beliefs about the distribution of types in the population are. Notice that condition (2) can not be binding, if investments are *substitutes*. Investments are *substitutes* if the marginal return from the first investment is greater than from the second investment:

$$\begin{aligned} S_{LH} - S_{LL} &> S_{HH} - S_{LH}, \\ \frac{2S_{LH}}{S_{LL} + S_{HH}} &> 1. \end{aligned}$$

Conversely, investments are complements if the marginal return from the second investment is larger:

$$\begin{aligned} S_{LH} - S_{LL} &< S_{HH} - S_{LH}, \\ \frac{S_{LL} + S_{HH}}{2S_{LH}} &> 1. \end{aligned}$$

In this case, when  $\delta$  is large it may be more efficient for *HL* pairs not to trade, and instead to delay trade until they meet a partner of the same type. For further reference, the requirement that there is a bargaining norm that implies trade in HL pairings regardless of the individual beliefs about the type distribution is summarized as the *trade condition*:

**Definition 1** *The discount rate  $\delta$  satisfies the trade condition if  $\delta < \frac{2S_{LH}}{S_{LL} + S_{HH}}$ .*

It shall be shown below that this is a necessary condition for the existence of a stochastically stable bargaining convention when investments are complements. By a *convention* we mean a pair  $\{I, \hat{x}_{LH}\}$ , with

the interpretation that each agent selects the investment  $I$  upon entering the market, the low type demands  $\hat{x}_{LH}$ , while the high type demands  $\hat{x}_{HL} = S_{LH} - \hat{x}_{LH}$ . To economize on writing out the full set of strategies and payoffs, the notion of a stable convention is defined as follows.

**Definition 2** *A convention  $\{H, \hat{x}_{LH}\}$  is stable if:*

1.  $(1 - \lambda)(S_{HH}/2 - \hat{x}_{LH}) + \lambda((S_{LH} - \hat{x}_{LH}) - \frac{S_{LL}}{2}) \geq c/(1 - 2\lambda)$ ,
2.  $S_{LH} - \hat{x}_{LH} \geq \delta \frac{(1-\lambda)}{(1-\delta\lambda)} S_{HH}/2$
3.  $\hat{x}_{LH} \geq \delta \frac{\lambda}{(1-\delta(1-\lambda))} S_{LL}/2$ .

The expected payoff of a person making a high investment assuming that trade is immediate and she meets a high type is  $(1 - \lambda)S_{HH}/2 + \lambda\hat{x}_{LH}$ , while the result of no investment is  $\lambda S_{HH}/2 + (1 - \lambda)\hat{x}_{LH}$ . If she meets a low type, the expected payoffs are  $(1 - \lambda)(S_{LH} - \hat{x}_{LH}) + \lambda S_{LL}/2$  if she invests high and  $\lambda(S_{LH} - \hat{x}_{LH}) + (1 - \lambda)S_{LL}/2$  if she invests low. Given the equilibrium fraction of high types in the market in any period is  $(1 - \lambda)$  a simple calculation yields condition 1. The second condition is the requirement that a person who is a high type prefers to trade with a low type, rather than wait until meeting a high type. The final condition requires the low type to prefer trading with a high type, rather than waiting until meeting a low type. This places a lower bound on  $\hat{x}_{LH}$ . It is a straightforward exercise to show that for every stable convention there is a subgame perfect Nash equilibrium yielding this outcome for the trading game outlined above. A stable convention,  $\{L, \hat{x}_{LH}\}$ , for low investment is defined in a similar fashion.

For much of the analysis the parameter  $\lambda$  is positive, but small. In the limit, when  $\lambda = 0$  then a sufficient condition for the existence of a stable convention with high investment is that it is efficient.

**Proposition 1** *Suppose it is strictly efficient for all agents to select high investment,  $S_{HH} - 2c > \max\{S_{LH} - c, S_{LL}\}$ , then for all  $\delta$  satisfying the trade condition a bargaining convention,  $\hat{x}_{LH}$ , exists such that  $\{H, \hat{x}_{LH}\}$  is stable for  $\lambda$  sufficiently small.*

This result demonstrates that when noise is small it is possible to support as an equilibrium high investment whenever it is efficient to do so. In contrast, the literature on the holdup problem assumes that the *ex post* division of the surplus is determined by the Nash bargaining solution, which in some cases induces sub-efficient investment. However the division implied by the Nash bargaining solution is only one among many subgame perfect equilibria of the game. In general, one is able to conclude that for this game there are a large number of subgame perfect equilibria, many of which induce efficient investment. The questions then is whether or not the efficient equilibria are stochastically stable.

## 4 Learning Dynamics

Consider now the kind of bargaining and investment conventions that are developed endogenously in a population of adaptive agents. Following Young (1993) and Tröger (2000) it is assumed that agents sample the previous periods trades to build an empirical distribution regarding the investment and bargaining

behavior of the other individuals in the population. Regarding the value of the outside option, agents believe that the distribution of low and high types in the economy is time stationary, a hypothesis that is consistent with assumption that agents base current action on past observations of the frequency of high types. It is also assumed that with a small probability they make mistakes in executing their optimal strategy given their beliefs regarding the play of the game described in section 2.

Our model consists of a single population of individuals who choose investment from  $\{h, l\}$  upon entering the population and afterwards every period have to choose their action from the space  $X(k)$ . This choice is based on beliefs about distribution of types and bargaining behavior of the other individuals in the population. Each period every individual independently takes a random sample of  $m$  individuals from the previous period. Let  $\hat{p}_t^i \in P = \{0, 1/m, 2/m, \dots, 1\}$  denote the fraction of individuals in this sample with  $T_{i,t-1} = H$ . Since the equal split occurs in all  $HH$  or  $LL$  pairings, only those observations where individuals select either  $x_{HL}$  or  $x_{LH}$  are useful for the estimation of bargaining behavior. These observations are used to update ones memory which is then used to estimate bargaining behavior of high and low types. The memory consists of at most  $m$  data points at any time, where it is assumed that the oldest data is dropped as new data is added. Observations in the memory are used to estimate empirical distribution functions  $\hat{F}_{HL}(\cdot)$  and  $\hat{F}_{LH}(\cdot)$  of bids of high types when matched with low types and vice versa. These distribution functions are taken from the finite set:

$$\mathcal{F} = \{F : X(k) \rightarrow \{0, 1/m, 2/m, \dots, 1\} \mid F(x) \text{ is increasing, } F(S_{LH}) = 1\}.$$

It will turn out to be convenient to denote by  $\mathcal{P}(z)$  the distribution function of point expectations  $z$ , i.e.  $\mathcal{P}(z)(x) = 0$  for  $x < z$  and  $\mathcal{P}(z)(x) = 1$  for  $x \geq z$ . When an agent leaves the market, her beliefs are passed on to the new agent entering the market to replace this agent. Beliefs in the first period are arbitrary.

The set of all possible beliefs of an agent is then given by  $B = P \times \mathcal{F}^2$ , where  $\hat{p}(\beta)$  and  $\hat{F}_{HL}(x_{HL}, \beta)$  denote respectively the proportion of high types and probability that  $x_{HL}$  or less is demanded by a high type given the belief  $\beta \in B$ . The expected payoff of an agent with type  $H$  or  $L$  choosing  $a \in X(k)$  under beliefs  $\beta \in B$ , is given recursively by:

$$\begin{aligned} U_L(a, \beta) &= \hat{p}(\beta) \left( \hat{F}_{HL}(S_{LH} - a, \beta) a + \delta \left( 1 - \hat{F}_{HL}(S_{LH} - a, \beta) \right) U_L(a, \beta) \right) + (1 - \hat{p}(\beta)) S_{LL}/2, \\ U_H(a, \beta) &= \hat{p}(\beta) S_{HH}/2 + (1 - \hat{p}(\beta)) \left( \hat{F}_{LH}(S_{LH} - a, \beta) a + \delta \left( 1 - \hat{F}_{LH}(S_{LH} - a, \beta) \right) U_H(a, \beta) \right). \end{aligned}$$

The time-line of the game with adaptive dynamics is summarized as follows:

1. At the beginning of the game beliefs are random, but when an individual leaves she is replaced by another agent with the same beliefs, say  $\beta$ .
2. Given beliefs  $\beta \in B$  the agent chooses to invest if:

$$\max_{(x_{LH}, x_{HL}) \in X(k)^2} (1 - \lambda) U_H(x_{HL}, \beta) + \lambda U_L(x_{LH}, \beta) - c \geq \max_{(x_{LH}, x_{HL}) \in X(k)^2} (1 - \lambda) U_L(x_{LH}, \beta) + \lambda U_H(x_{HL}, \beta).$$

Then she draws her type, which is equal to her investment with probability  $1 - \lambda$ .



Each period the following steps are repeated until exit occurs:

1. At the beginning of every period  $t$  the individual randomly samples the types of  $m$  individuals from the previous period. This is used to update beliefs  $b_t^i \in B$ .
2. With probability  $\varepsilon > 0$  the individual selects an action randomly from  $X(k)$ , under the uniform distribution. This noise process is *i.i.d.* between individuals and periods. With probability  $1 - \varepsilon$  the individual chooses  $a_t^i \in X(k)$  to maximize  $U_{T^i}(a_t^i, b_t^i)$ , given her type  $T^i \in \{L, H\}$  and beliefs  $b_t^i \in B$ . When indifferent over demands she chooses the smallest demand. The agent's strategy is uniquely defined by her beliefs. Hence, we write  $a_t^i = \alpha(T^i, b_t^i)$ .
3. Agents are randomly paired, and their payoffs are determined according to the actions chosen at stage 2.

Given that an agent's action is completely characterized by her beliefs  $b_t^i \in B$ , and type  $T^i \in \{H, L\}^2$ , the state at time  $t$  is characterized by a distribution over beliefs and types, and the state space is therefore finite and given by:

$$(3) \quad \mathcal{S} = \{s \in [0, 1]^{|C|} \mid \sum_{c \in C} s_c = 1, \quad ns_c \in \mathbf{N}_0 \quad \forall c \in C\},$$

where  $C = \{H, L\} \times B$ . The learning process described above defines a time homogeneous Markov process  $\{\sigma_t\}_{t=0}^\infty$  on the state space  $\mathcal{S}$ . Although, even for  $\varepsilon > 0$ , the transition matrix is not positive, the following lemma shows that the process is irreducible and aperiodic.

**Lemma 1** *For  $\varepsilon > 0$  the Markov process  $\{\sigma_t\}_{t=0}^\infty$  as defined above is irreducible and aperiodic.*

Hence, for  $\varepsilon > 0$  there exists a unique limit distribution  $\pi^*(\varepsilon)$  over  $\mathcal{S}$ , where  $\pi_s^*(\varepsilon)$  denotes the probability of state  $s$ . Following a standard approach in evolutionary game theory we consider the limit distribution for small values of  $\varepsilon$  and in particular characterize the states whose weight in the limit distribution stays positive as the mutation probability  $\varepsilon$  goes to 0. Such states are called stochastically stable:

**Definition 3** *A state  $s \in \mathcal{S}$  is called stochastically stable if  $\lim_{\varepsilon \rightarrow 0} \pi_s^*(\varepsilon) > 0$ . We say that a set is stochastically stable if all his elements are stochastically stable.*

The reason why this concept is of interest is that for small  $\varepsilon$  the process spends almost all the time in stochastically stable sets. Hence, characterizing the stochastically stable outcome means characterizing the long run properties of the evolutionary process. To identify stochastically stable states it is necessary to first identify the minimal absorbing sets of the process for  $\varepsilon = 0$ . It is well known that the set of stochastically stable states is a subset of the union of these so called limit sets. Formally, a limit set is defined as follows:

**Definition 4** *A set  $\Omega \subseteq \mathcal{S}$  is called a limit set of the process if for  $\varepsilon = 0$  the following statements hold:*

$$\begin{aligned} \forall s \in \Omega \quad \mathbb{P}(\sigma_{t+1} \in \Omega \mid \sigma_t = s) &= 1 \\ \forall s, \tilde{s} \in \Omega \quad \exists z > 0 \text{ s.t. } \mathbb{P}(\sigma_{t+z} = \tilde{s} \mid \sigma_t = s) &> 0. \end{aligned}$$

---

<sup>2</sup>We look at the process after all incoming agents have made their investment decisions, but before they are paired and therefore the type of all agents is determined.

In the following section we will characterize the stochastically stable sets and discuss the implied investment and bargaining conventions. This will allow us to highlight the different implications for investment such an evolutionary perspective has compared to assuming that the Nash bargaining solution is used.

## 5 Stochastically Stable Conventions

The question we address is the emergence of a unique, efficient and stable bargaining convention in which all individuals follow the same investment strategy, and have the same expectations regarding how to divide the gains from trade. This is formally defined by:

**Definition 5** *A state  $s$  induces the bargaining convention  $x_{LH}$  if all individuals have beliefs  $\beta \in B$  that place probability one on the demand by their partner being  $x_{LH}$  or  $S_{LH} - x_{LH}$ , depending upon their type in  $HL$  matches.<sup>3</sup>*

Therefore, we shall say that a bargaining convention does not exist at a state  $s$  if there is heterogeneity in the beliefs of the agents regarding the terms of trade between high and low types. Let us now consider the constraints that the outside options place upon feasible bargaining conventions.

A necessary condition for a convention is that the terms of trade between high and low types result in outcomes that are better than their respective the outside options. Consider an agent with beliefs  $\beta$  then, under the assumption of stationary beliefs, by simply waiting for a partner with the same type she can guarantee an expected payoff of  $\frac{\hat{p}(\beta)}{1-\delta(1-\hat{p}(\beta))} \frac{S_{HH}}{2}$  if she is of type  $H$  and  $\frac{1-\hat{p}(\beta)}{1-\delta\hat{p}(\beta)} \frac{S_{LL}}{2}$  if she is of type  $L$ . We say that a bargaining convention is compatible with  $\hat{p}$  and  $\delta$  if both parties are better off than their respective expected outside option, that is  $x_{LH} \in [\underline{x}_{LH}(\hat{p}), \bar{x}_{LH}(\hat{p})]$ , where  $\underline{x}_{LH}(\hat{p}) \in X(k)$  such that

$$\underline{x}_{LH}(\hat{p}) - \alpha < \frac{\delta(1-\hat{p})}{1-\delta\hat{p}} \frac{S_{LL}}{2} \leq \underline{x}_{LH}(\hat{p})$$

and  $\bar{x}_{LH}(\hat{p}) \in X$  such that

$$\bar{x}_{LH}(\hat{p}) \leq S_{LH} - \frac{\delta\hat{p}}{1-\delta(1-\hat{p})} \frac{S_{HH}}{2} < \bar{x}_{LH}(\hat{p}) + \alpha,$$

where  $\alpha = \alpha_{LH} = \alpha_{HL}$  is the minimum unit of account for dividing the surplus, as defined in the game form of section 2. Denote the set of all bargaining conventions which are compatible with all  $\hat{p} \in [0, 1]$  for a certain discount factor by  $\mathcal{C}(\delta) = [\underline{x}_{LH}(0), \bar{x}_{LH}(1)]$ . Notice, that  $\mathcal{C}(\delta) \neq \emptyset$  for sufficiently small  $\alpha$  if and only if  $\delta < \frac{2S_{LH}}{S_{LL}+S_{HH}}$  holds. Hence, the trade condition is a sufficient condition for  $\mathcal{C}(\delta) \neq \emptyset$ , and is also a necessary condition in the case of complementary investments.

Let us now characterize the limit sets in this framework. Once a bargaining convention  $x_{LH}$ , which is compatible with  $\delta$ , is reached, in the absence of mutations all low types always demand  $x_{LH}$  against high types and high types always demand  $S_{LH} - x_{LH}$  against low types. Hence beliefs can never change once such a state has been reached. If beliefs are heterogeneous there is always a positive probability that all agents observe identical samples and beliefs become homogeneous and compatible. However, also after a

<sup>3</sup>Formally  $\hat{F}_{HL}(\beta) = \mathcal{P}(S_{LH} - x_{LH})$ , and  $\hat{F}_{LH}(\beta) = \mathcal{P}(x_{LH})$ .

bargaining convention has been reached, the distribution of agent types may change between two periods, even if the investment behavior is constant. The randomness of the outcome from investment implies that all distributions of  $H$  and  $L$  types are possible. Hence if there is a bargaining convention that is not compatible with all  $p \in P$ , eventually it will be disrupted. This suggests that the limit sets correspond to conventions in  $\mathcal{C}(\delta)$ , when it is not empty.

If the trade condition does not hold then  $\mathcal{C}(\delta) = \emptyset$ , and the outside option of waiting for an equal type always becomes binding for some  $\hat{p}$ . In this case bids never settle down at a compatible convention and there occur fluctuations driven by the fluctuations in the  $\hat{p}^i$ . In the following lemma we show that in such a scenario the set of possible bids is given by all demands which lie just above the outside option for some  $\hat{p} \in P$  and the best responses to that. The set of all these bids is given by

$$\begin{aligned}\mathcal{B}_{LH}(\delta) &= \{x_{LH} \in X | \exists p \in P \text{ s.t. } x_{LH} = \underline{x}_{LH}(p) \text{ or } \exists p \in P \text{ s.t. } x_{LH} = \bar{x}_{LH}(p)\} \\ \mathcal{B}_{HL}(\delta) &= \{x_{HL} \in X | \exists p \in P \text{ s.t. } x_{HL} = S_{LH} - \bar{x}_{LH}(p) \\ &\quad \text{or } \exists p \in P \text{ s.t. } x_{HL} = S_{LH} - \underline{x}_{LH}(p)\}.\end{aligned}$$

The larger  $m$  is the larger these sets are and for sufficiently large  $m$  we simply have  $\mathcal{B}_{HL}(\delta) = X \cap [S_{LH} - \frac{\delta S_{LL}}{2}, \frac{\delta S_{HH}}{2}]$  and  $\mathcal{B}_{LH} = X \cap [S_{LH} - \frac{\delta S_{HH}}{2}, \frac{\delta S_{LL}}{2}]$ .

**Lemma 2** *Characterization of the limit sets for sufficiently large  $k$ :*

a) *Suppose the trade condition holds then for each  $x_{LH} \in \mathcal{C}(\delta)$  there exists a limit set  $\Omega(x_{LH})$  consisting of all  $s \in \mathcal{S}$  such that  $s_\zeta > 0$  only if  $\zeta = (T, \beta)$  for some  $T \in \{H, L\}$  and some  $\beta$  such that  $\hat{F}_{HL}(\cdot, \beta) = \mathcal{P}(S_{LH} - x_{LH})$  and  $\hat{F}_{LH}(\cdot, \beta) = \mathcal{P}(x_{LH})$ .*

b) *Suppose investments are complements and the trade condition does not hold, then there exists a single limit set  $\mathcal{L}$ . For all states  $s \in \mathcal{L}$  we have  $s_\zeta > 0$  only if  $\zeta = (T, \beta)$  for some  $T \in \{H, L\}$  and some  $\beta$  such that  $\text{supp}(\hat{F}_{HL}(\cdot, \beta)) \in \mathcal{B}_{HL}(\delta)$ ,  $\text{supp}(\hat{F}_{LH}(\cdot, \beta)) \in \mathcal{B}_{LH}(\delta)$ .*

Hence, as long as  $\delta$  is not too large or investments are not too complementary ( $S_{LH}$  is small), in the long run the process generates a bargaining convention which is followed by all individuals in the population. The bargaining strategy is uniform across the population, there is no disagreement and bargaining is always efficient. It should also be pointed out that conventions only fail to exist for large  $\delta$  if investments are complements ( $S_{LH} < \frac{1}{2}(S_{HH} + S_{LL})$ ). In that case the trade condition may not be satisfied for some  $\delta \in [0, 1)$ , and hence agents may choose not to trade in  $HL$  pairings, making it impossible for a convention to evolve in this case.

The analysis in the subsequent sections uses the Radius - Modified Coradius criterion, recently developed in Ellison (2000), to provide sufficient conditions for a limit set to be stochastically stable. A proof using the Freidlin-Wentzell Freidlin and Wentzell (1984) technique – which has been used in the analysis of most models of a similar kind – would also be feasible but slightly more complicated. The arguments used to determine the stochastically stable convention are similar to those used in Young (1993) but it turns out that the fact that individuals have the option to wait for a different partner has significant and interesting implications for the long run bargaining behavior.

## 6 The Case of Substitutes

Consider first the case in which investments are substitutes, namely  $S_{LH} - S_{LL} \geq S_{HH} - S_{LH}$ . In this case, the gains from having one person invest are greater than having a second person invest.

The fact that individuals cannot perfectly determine the productivity of their investments has two important implications. First, every period there is a strictly positive probability of both types existing in the market, and thus there are always with positive probability  $HL$  trades occurring in the market which can be used to update the beliefs of individuals. Second, regardless of the investment decisions of individuals, any distribution of types has strictly positive probability. On the other hand, transitions between bargaining conventions have to be triggered by (in general multiple simultaneous) mutations. Hence, for small mutation probabilities bargaining conventions adjust more slowly, and are more stable than the realized distribution of types. This implies that the stochastically stable bargaining convention is *independent* of the long run investment behavior and therefore also independent of investment costs  $c$ . The next proposition provides a rigorous proof of this fact and derives the properties of the bargaining convention that arises in the long run.

**Proposition 2** *For sufficiently large  $m, n$  the limit of the stochastically stable sets of the process  $\{\sigma_t\}$  for  $k \rightarrow \infty$  can be characterized in the case of substitutes as follows:*

(a) *When  $S_{HH} - \frac{\delta}{2}(S_{HH} - S_{LL}) \geq S_{LH} \geq (S_{HH} + S_{LL})/2$ , every stochastically stable state induces the bargaining convention*

$$\hat{x}_{LH}^s = \frac{S_{LH}}{2} - \frac{\delta}{2(2-\delta)}(S_{LH} - S_{LL}).$$

(b) *When  $S_{HH} \geq S_{LH} \geq S_{HH} - \frac{\delta}{2}(S_{HH} - S_{LL})$  every stochastically stable state induces the bargaining convention*

$$\hat{x}_{LH}^s = \frac{S_{LH}}{2} - \frac{\delta}{4}(S_{HH} - S_{LL}).$$

We will refer to the bargaining convention induced by all stochastically stable states as the stochastically stable bargaining convention. In the absence of outside options ( $\delta = 0$ ) the equal split rule is the unique, stochastically stable bargaining convention, regardless of the investment levels. Possibly the more surprising result is the effect of the outside options on the bargaining convention. Notice, that in this model the outside option is introduced only as a constraint on the set of possible bargaining agreements, and hence one might expect the outside option principle to apply (see Binmore, Rubinstein, and Wolinsky (1986)). In that case if  $x_{LH} > \delta S_{LL}/2$  and  $S_{LH} - x_{LH} > \delta S_{HH}/2$ , then  $x_{LH}$  should not depend on either  $S_{HH}$  or  $S_{LL}$ , yet we find that that for all  $\delta > 0$  the stochastically stable bargaining convention depends upon at least one of the outside options, and that the low types share is always strictly increasing in  $S_{LL}$ , a result that is consistent with Binmore, Proulx, and Samuelson (1995) who report results from a bargaining game with drift.

Though the outside option affects the outcome of bargains, the level of long run investment does not. This greatly simplifies the analysis of investment behavior in the long run. We can determine investment behavior as a function of the bargaining convention and then insert the stochastically stable bargaining

convention. For a given bargaining convention  $\hat{x}_{LH}$  investment is optimal iff

$$\begin{aligned} & (1 - \lambda)(\hat{p}\frac{S_{HH}}{2} + (1 - \hat{p})(S_{LH} - \hat{x}_{LH}) + \lambda(\hat{p}x_{LH} + (1 - \hat{p})\frac{S_{LL}}{2}) - c \\ & \geq (1 - \lambda)(\hat{p}\hat{x}_{LH} + (1 - \hat{p})\frac{S_{LL}}{2}) + \lambda(\hat{p}\frac{S_{HH}}{2} + (1 - \hat{p})(S_{LH} - \hat{x}_{LH})). \end{aligned}$$

Taking into account that  $(S_{HH} + S_{LL}/2 - S_{LH} < 0$  this gives the following condition for high investment to be optimal:

$$(4) \quad \hat{p} \leq p^*(\hat{x}_{LH}; \lambda) := \frac{S_{LH} - \hat{x}_{LH} - S_{LL}/2 - c/(1 - 2\lambda)}{S_{LH} - S_{HH}/2 - S_{LL}/2}.$$

To analyze the dynamics of investment for a given bargaining convention we consider the evolution of type distributions over time. Given a current distribution of types the distribution of types in the following period in general depends on the outcome of the stochastic sampling procedure for all agents, which gives the beliefs  $\hat{p}(b_t^i)$  and therefore influences the investment decisions, and the actual realization of types given the investment decision. This can be described by a Markov process  $\{\tilde{\sigma}_t\}_{t=0}^{\infty}$  on the state space  $\tilde{S} = \{0, 1/n, 2/n, \dots, 1\}$ . For  $\lambda > 0$  the process is irreducible and aperiodic. The unique limit distribution is denoted by  $\tilde{\pi}^*(\lambda)$ . The following lemma characterizes the limit distribution for small values of  $\lambda$ .

**Lemma 3** *When investments are substitutes ( $S_{LH} > \frac{1}{2}(S_{HH} + S_{LL})$ ), then given a bargaining convention  $\hat{x}_{LH}$ , the long run distribution of types for sufficiently small  $\lambda$  can be characterized as follows:*

(a)  $p^*(\hat{x}_{LH}; 0) \leq 0$ : no individual ever invests and  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_0^*(\lambda) = 1$ .

(b)  $p^*(\hat{x}_{LH}; 0) > 1$ : all individuals always invest and  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_1^*(\lambda) = 1$ .

(c)  $p^*(\hat{x}_{LH}; 0) \in (0, 1)$ :  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_1^*(\lambda) = \lim_{\lambda \rightarrow 0} \tilde{\pi}_0^*(\lambda) = 0.5$ .

In case (a) we say that  $\hat{x}_{LH}$  induces a no-investment convention, in (b)  $\hat{x}_{LH}$  induces a full investment convention, and in case (c) we say that  $\hat{x}_{LH}$  induces cyclical investment. By cyclical investment we mean that in one period everybody invests, and in the next period nobody invests. What is happening is that when all individuals invest, it is optimal not to invest, and verso. Note that for substitutes in cases where  $p^*(\hat{x}_{LH}; 0) > 1$  the action  $H$  is dominant at the investment stage for small  $\lambda$  and all heterogeneity in types is created by deviations of the actual type from investment. Therefore, it is easy to see that a bargaining convention  $\hat{x}_{LH}$  induces an investment convention if and only if there is a  $\lambda^* > 0$  such that the convention  $\{H, \hat{x}_{LH}\}$  is stable for all  $\lambda < \lambda^*$ . It should also be pointed out here that even if we assume that  $\lambda$  is small it is still assumed to be of an order of magnitude larger than the mutation probability  $\epsilon$  which means that the transition between bargaining conventions is always assumed to be much slower than the transition between investment patterns. Using this result it is straight forward to describe the investment behavior which is induced by the stochastically stable bargaining conventions. Inserting the stochastically stable bargaining convention  $\hat{x}_{LH}^s$  into  $p^*$  and applying lemma 3 gives the following proposition.

**Proposition 3** Assume that  $m, n$  and  $k$  are sufficiently large.

(a) For  $\frac{1}{2}(S_{HH} + S_{LL}) < S_{LH} \leq S_{HH} - \frac{\delta}{2}(S_{HH} - S_{LL})$  the stochastically stable bargaining convention induces full-investment for  $c < c^1$ , no-investment for  $c > c^2$  and cyclical investment for  $c \in [c^1, c^2]$ , where

$$\begin{aligned} c^1 &= \frac{1}{2(2-\delta)}(\delta(S_{LH} - S_{LL}) + (2-\delta)(S_{HH} - S_{LH})) \\ c^2 &= \frac{1}{2-\delta}(S_{LH} - S_{LL}). \end{aligned}$$

(b) For  $S_{HH} - \frac{\delta}{2}(S_{HH} - S_{LL}) < S_{LH} \leq S_{HH}$  the stochastically stable bargaining convention induces full-investment for  $c < c^3$  and cyclical investment for  $c \geq c^3$ , where

$$c^3 = \frac{1}{4}(\delta(S_{HH} - S_{LL}) + 2(S_{HH} - S_{LH})).$$

Notice that when  $\delta = 0$ , then  $c_1 = (S_{HH} - S_{LH})/2$ , but in this case of substitutes it is efficient for both parties to invest whenever  $c < (S_{HH} - S_{LH})$ . Therefore we obtain under-investment in some cases.

In case (a), the gain from investing at the bargaining convention is:

$$\begin{aligned} S_{HH}/2 - x_{LH} &= \frac{(S_{HH} - S_{LH})}{2} + \frac{\delta}{2(2-\delta)}(S_{LH} - S_{LL}), \\ &\geq \frac{(S_{HH} - S_{LH})}{2}. \end{aligned}$$

Therefore, the outside option always increases the gains from investing, regardless of whether it is binding at the equilibrium. However, for case (a) under the trade condition, it never increases incentives to the point that the gains from investing are equal to the full marginal gains, given by  $(S_{HH} - S_{LH})$ . On the other hand, if investments are strong substitutes and the gains from the second investment are very small (case (b) above) the stochastically stable convention indeed induces full investment whenever this is efficient. This is formalized in the following corollary.

**Corollary 1** For  $S_{HH} - \frac{\delta}{2}(S_{HH} - S_{LL}) < S_{LH} \leq S_{HH}$  the stochastically stable bargaining convention induces full-investment for all values of  $c$  where full investment is efficient.

**Proof.** It is straight-forward to check that  $c^3 \leq S_{HH} - S_{LH}$  under these assumptions. ■

As pointed out above, our results about the stochastically stable bargaining conventions show that the outside option acts rather as a threat point in the allocation of the joint surplus. This might raise the question whether the efficiency result of corollary 1 is a simple implication of the difference in threat point payoffs of the two types. To address this question let us denote by  $\hat{x}_{LH}^N$  the allocation consistent with the Nash bargaining solution between a high and a low type where both have beliefs  $\hat{p} = 1$  and the expected payoffs in the following period are treated as a threat point. This allocation has to satisfy

$$\hat{x}_{LH}^N = \delta \hat{x}_{LH}^N + \frac{1}{2} \left( S_{LH} - \delta \hat{x}_{LH}^N - \delta \frac{S_{HH}}{2} \right),$$

and therefore we get

$$(5) \quad \hat{x}_{LH}^N = \frac{S_{LH}}{2} - \frac{\delta(S_{HH} - S_{LH})}{2(2-\delta)}.$$

Comparing this expression with the stochastically stable bargaining conventions from proposition 2, simple calculations show that under our assumption of  $S_{LH} > (S_{HH} + S_{LL})/2$  we always have  $\hat{x}_{LH}^N > x_{LH}^s$ . Accordingly, the investment incentives in a population of investors under the stochastically stable bargaining norm are not only larger than under the equal split rule but also larger than under Nash bargaining with the outside options as threat points. To understand this result intuitively we have to realize that the long run stability of the bargaining norms are determined by their resistance to change in scenarios where deviations from the norm have the highest chance of altering the norm. Bargaining norms are the easiest destabilized in scenarios with low investment in the population since the amount an agent risks when following deviators from the norm is the smallest under this investment pattern. If investments are substitutes a high type has a lot of bargaining power in an environment of low types and hence the stochastically stable bargaining norm gives a larger part of the surplus to the high types than they would get if the norm had been evolved in a population of mostly high types. Hence, the stochastically stable convention allocates more to the high types than the Nash bargaining solution in an environment of high types would. Since the stochastically stable bargaining norm although developed in low investment scenarios is adhered to even if in the long run everyone invests, it facilitates the development of full investment norms.

This discussion implies that the evolutionary approach facilitates the development of full investment. In the following corollary we compare the stochastically stable outcome to the notion of a stable convention under the equal split rule and the Nash bargaining solution:

**Corollary 2** •

(a) *If  $c$  satisfies*

$$\frac{1}{2}(S_{HH} - S_{LH}) < c < \frac{\delta}{4}(S_{HH} - S_{LL}) + \frac{1}{2}(S_{HH} - S_{LH}),$$

*then for  $\lambda$  sufficiently small the stochastically stable convention induces full investment, but  $\{H, S_{LH}/2\}$  is not a stable convention.*

(b) *If  $c$  satisfies*

$$\frac{1}{2-\delta}(S_{HH} - S_{LH}) < c < \frac{\delta}{4}(S_{HH} - S_{LL}) + \frac{1}{2}(S_{HH} - S_{LH}),$$

*then for  $\lambda$  sufficiently small the stochastically stable convention induces full investment, but  $\{H, \hat{x}_{LH}^N\}$  is not a stable convention.*

Two remarks concerning this corollary are in order. First, it is easy to see that the ranges of  $c$  given in parts (a) and (b) are both non-empty if investments are substitutes and  $\delta > 0$ . Second, part (b) shows that if the allocation of surplus is determined by the Nash bargaining solution even with the outside option as threat point there always remains a hold-up region, i.e. a range of parameters where full investment is efficient but  $\{H, \hat{x}_{LH}^N\}$  is not a stable convention. On the other hand, this is not the case for the stochastically stable bargaining convention. This result illustrates that in the case of substitutes, endogenously determined bargaining conventions yield a larger set of parameter values with high investment than cooperative solutions. Consider now the case of complements.

## 7 The Case of Complements

Consider the case of complementary investments, where  $(S_{HH} + S_{LL})/2 \geq S_{LH} \geq S_{LL}$ . This implies that  $S_{HH} - S_{LH} \geq S_{LH} - S_{LL}$ , and hence the marginal gain from investment is higher after one person has invested. If the trade condition is not satisfied, then for  $S_{LH}$  sufficiently close to  $S_{LL}$ , either the high types prefer to wait for a high type rather than trade with a low type if  $\hat{p}$ , the fraction of higher types, is sufficiently high, or the low types prefer to wait for a low type rather than trade with a high type if  $\hat{p}$ , the fraction of higher types, is sufficiently low. When this occurs, a stable norm of behavior does not evolve, as shown in the following proposition.

**Proposition 4** *Suppose that investments are complements and the trade condition does not hold, that is  $\mathcal{C}(\delta) = \emptyset$ , then for sufficiently large  $m$ ,  $n$  and  $k$  the unique stochastically stable sets of the process  $\{\sigma_t\}$  is  $\mathcal{L}$ , as defined in lemma 2, and there exist no stochastically stable bargaining conventions.*

**Proof.** This follows immediately from Lemma 2, which shows that the set  $\mathcal{L}$  is the only limit set under these conditions. ■

Therefore, a necessary condition for the evolution of a bargaining convention is  $\mathcal{C}(\delta) \neq \emptyset$ , where regardless of the fraction of high types expected in the market, there exist bargaining conventions such that both high and low types prefer to trade rather than wait. In this case, there is a unique stochastically stable bargaining convention which again does not depend upon beliefs regarding the fraction of high types in the market.

**Proposition 5** *Suppose the trade condition holds, then for sufficiently large  $m$ ,  $n$  the limit of the stochastically stable sets of the process  $\{\sigma_t\}$  for  $k \rightarrow \infty$  can be characterized as follows:*

(a) For  $S_{LL} \leq S_{LH} \leq \frac{\delta}{2}((2 - \delta)S_{HH} + S_{LL})$  the stochastically stable bargaining convention is

$$\hat{x}_{LH}^s = S_{LH} - \delta \frac{S_{HH}}{2}.$$

(b) For  $\frac{\delta}{2}((2 - \delta)S_{HH} + S_{LL}) < S_{LH} \leq (S_{HH} + S_{LL})/2$  the stochastically stable bargaining convention is

$$\hat{x}_{LH}^s = \frac{S_{LH}}{2} - \frac{\delta}{2(2 - \delta)}(S_{LH} - S_{LL}).$$

Case (a) occurs when the outside option for the high type is binding for  $\hat{p} = 1$ . A necessary and sufficient condition for this case to apply is:

$$1 - \sqrt{1 - 2 \left( \frac{S_{LH} - S_{LL}}{S_{HH}} \right)} \leq \delta \leq \frac{2S_{LH}}{S_{HH} + S_{LL}}.$$

If the discount factor is too high, then individuals do not wish to enter into  $HL$  trade. Conversely, if the discount factor is low, then the outside option for the high type is not binding. However, as case (b) illustrates, one of the implications of stochastic stability criteria is that, as in the case of substitutes, the existence of an outside options *always* increases the payoff for the high type relative to the equal division solution. On the other hand, it can be easily verified that the Nash bargaining solution with the outside



option as threat point, given by (5), gives a smaller allocation of the surplus to low types compared to the stochastically stable convention and therefore provides higher investment incentives.

With complements investment incentives are larger for  $\hat{p} = 1$  than for  $\hat{p} = 0$ . This implies that if no investment is optimal at  $\hat{p} = 1$ , no individual will invest any more, once the bargaining convention  $\hat{x}_{LH}$  has been established – a no-investment convention is induced. On the other hand, if investment is optimal at  $\hat{p} = 0$ , everyone invests under the stochastically stable bargaining convention – a full investment convention is induced. However, if investment is optimal for  $\hat{p} = 1$  and no investment is optimal for  $\hat{p} = 0$ , both the homogeneous state corresponding to full investment and the homogeneous state corresponding to no investment are locally stable states in the sense that the process never leaves each of these states as long as high investment always implies high types and low investment always implies low types. In such a scenario the threshold  $p^*$  defined in (4) is in  $(0, 1)$  and investment is optimal if and only if  $\hat{p} \geq p^*(\hat{x}_{LH})$ .

Investment effects are however assumed to be stochastic, and therefore there is always a positive probability that the process wanders from a non-investment to a full investment state and vice versa. As in the case of substitutes, for a given bargaining convention the evolution of the distribution of high and low types is described by a Markov process  $\{\tilde{\sigma}_t\}_{t=0}^\infty$  on the state space  $\tilde{S}$ . Denoting again by  $\tilde{\pi}^*$  the unique limit distribution we get the following lemma:

**Lemma 4** *Assume that  $0 < \lambda < 0.5$  and  $m$  and  $n$  are and sufficiently large and a bargaining convention  $x_{LH}$  is given. Then, for  $p^*(x_{LH}, \lambda) > (<)0.5$  we have  $\sum_{i < n/2} \tilde{\pi}_{i/n}^*(\lambda) > (<) \sum_{i > n/2} \tilde{\pi}_{i/n}^*(\lambda)$ . Furthermore, we have  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_0^* = 1$  if  $p^*(x_{LH}, 0) > 0.5$  and  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_1^* = 1$  if  $p^*(x_{LH}, 0) < 0.5$ .*

According to this lemma  $p^*(x_{LH}, 0) < 0.5$  implies that in the long run the probability to have a majority of high types is larger than the probability to have a majority of low types and as  $\lambda$  goes to zero the probability to see only high types goes to one. Note that in the case of complements the investment stage has the structure of a coordination game and hence this lemma basically rephrases well known results by Kandori, Mailath, and Rob (1993). We say that a no-investment convention is induced if the threshold  $p^*(x_{LH}, 0)$ , is larger than 0.5 and that a full investment convention is induced if this inequality holds the other way round. Using this we get the following characterization of the investment conventions induced by stochastically stable bargaining conventions.

**Proposition 6** *Assume that  $m, n$  and  $k$  are sufficiently large, the trade condition holds, and investments are complements, then the stochastically stable bargaining convention induces full investment if  $c < c^A(S_{LH}, \delta)$  and no-investment for  $c > c^A(S_{LH}, \delta)$ , where*

$$(6) \quad c^A(S_{LH}, \delta) = \begin{cases} \frac{1}{4}(S_{HH} - S_{LL}) + \frac{1}{2}(\delta S_{HH} - S_{LH}) & \text{if } S_{LH} \leq \frac{\delta}{2}((2 - \delta)S_{HH} + S_{LL}), \\ \frac{1}{4}(S_{HH} - S_{LL}) + \frac{\delta}{2(2 - \delta)}(S_{LH} - S_{LL}) & \text{if not.} \end{cases}$$

It follows from the coordination game structure of the investment stage that a bargaining norm  $\hat{x}_{LH}$  does not necessarily induce a high investment convention even if  $\{H, S_{LH}/2\}$  is a stable convention. An interesting implication of this insight, especially when compared to the case of substitutes, is that the set of parameters for which a full investment convention is stable under the equal split rule is *larger* than the

set of parameter values for which high investment is part of a stochastically stable equilibrium. To see this, notice that if  $\lambda = \delta = 0$ , then  $\{H, S_{LH}/2\}$  is a stable convention if and only if:

$$\frac{1}{2}(S_{HH} - S_{LH}) \geq c.$$

Therefore the following result is immediate.

**Corollary 3** *When  $\delta$  and  $\lambda$  are sufficiently small, then if  $c$  satisfies:*

$$\frac{1}{4}(S_{HH} - S_{LL}) < c < \frac{1}{2}(S_{HH} - S_{LH}),$$

*$\{H, \frac{S_{LH}}{2}\}$  is a stable convention, but there is no stochastically stable convention with full investment.*

Clearly, if the Nash bargaining solution  $\hat{x}_{LH}^N$  would be considered instead of the equal split rule the region with a stable high investment convention but no stochastically stable convention with high investment would be even larger. Overall these results illustrate that in the case of complements, stochastic stability implies that the holdup problem is even more severe than under the assumption of cooperative bargaining solutions.

## 8 Discussion

In this model we have assumed that, when bargaining over the joint surplus, each individual makes her bid contingent on the type of the two partners (i.e. on their contribution to the joint surplus) but not explicitly on the investment. If one would assume that investment itself is observable as well and taken into account by the bidders, a bidding strategy would have to specify a bid for each combination of investment and type of the two players. In ?) it has been shown that in a framework, where investment and type coincide with certainty, bargaining conventions are never established and the set of values of  $c$  and  $S_{LH}$  where long run investment conventions evolve is small compared to the set where investment is efficient. The problem is that in order for a convention of behavior to develop, it must be observed in the long run. When investments are observable, and high investment is the desired equilibrium, then  $LH$  trades would not be observed. Consequently, beliefs regarding the appropriate division in this cases tend to drift around, and an efficient convention cannot be sustained. This implies that an increase in the amount of observable information would yield a decrease in the long-run efficiency. It should be pointed out, that the fact that we consider two-sided investment is essential here. Tröger (2000) and Ellingsen and Robles (2000) assume deterministic investment effects in their analyses of the one-sided investment case and there the drift of beliefs is the driving force behind their efficiency results.

The case of stochastic investment does ensure the evolution of a bargaining convention whenever the trade condition is satisfied, in other words, as long as individuals find it in their interest to always trade, and there are always a significant number of  $HL$  trades occurring. In contrast to the earlier results for the one-sided investment case we find that for  $\delta = 0$  the bargaining convention ignores prior investment. This is a key ingredient for holdup to occur, as illustrated in Grout (1984) and Grossman and Hart (1986) who assume the terms of trade are determined by the Nash bargaining solution.

The amount of long run investment in our model depends crucially upon whether investments are substitutes or complements. In the case of substitutes, the level of investment is high for a wider range of parameter values as compared to the standard holdup model based upon either the equal split rule or the Nash bargaining solution. Conversely, with complements the situation is worse, and therefore we conclude that when learning is incorporated into investment behavior one obtains quite different results compared to the static model. Moreover, whether or not learning makes the situation better or worse is sensitive to the degree of complementarity between investments.

The distinction between our results and those of Tröger (2000) and Ellingsen and Robles (2000) have analogies in the contract theory literature. In the case of one sided investment when trade is always efficient then the results of Grossman and Hart (1986) and Hart and Moore (1990) illustrate a number of mechanisms can be used to achieve efficient investment, including the use of fixed price contracts and the appropriate allocation of property rights.<sup>4</sup> Our model corresponds to the case of cooperative investment, that is both parties make investments that affect the other person's gain.<sup>5</sup> In this case, as Che and Hausch (1999) and Hart and Moore (1999) observe, there does not exist any renegotiation proof mechanism that can implement the first best. Moreover, even if the mechanism entails inefficient punishments, since  $S_{LH} = S_{HL}$  at least for  $\delta = 0$  there is no way to screen between the  $H$  and  $L$  types *ex post*, hence any mechanism uniquely implementing high investment would entail some social loss.

A new element of the current paper, relative to the previous work, is that we consider the effect of potential competition. One reason for exploring this effect is that earlier work by MacLeod and Malcolmson (1993), and more recently Felli and Roberts (2000) and Cole, Mailath, and Postlewaite (2000) have shown that potential competition, even if imperfect, may completely solve the holdup problem. In this paper we find that imperfect competition acts more like a threat point in the Nash Bargaining sense, than as outside options. This is found to always increase the incentives for efficient investment, and in the case where investments are strong substitutes this leads to efficient investment.

Finally, it is the case that in this model there exist efficient equilibria, all of which have the property that the terms of trade are sensitive to sunk costs. Carmichael and MacLeod (1999) have shown that when there is sufficient diversity in payoffs the efficient rule is unique. Thus *efficient* fair division rule should incorporate the effect of sunk costs, which may explain why sunk costs may matter in decision making, as first observed by Thaler (1980). The results of this paper suggest that from the perspective of a simple two party investment problem, there appears to be a tension between sharing rules that are *ex ante* efficient, and those that are stochastically stable. Hence it is still an open question as to which precise process leads individuals to be sensitive to sunk costs. On the positive side, these results suggest that there are limits to development of efficient norms of behavior which are sensitive to the degree of complementarity between investments, and which may explain why efficiency is enhanced through the explicit introduction of social

---

<sup>4</sup>Also see the recent work of Robles (2001) who explores the evolution of contracts in the one-sided investment case.

<sup>5</sup>In the case of a buy-seller relationship, the self-investment case corresponds to the buyer making investments that affect the utility from consuming the good, while the seller's investment affects the cost of production. In this case a simple fixed price contract can achieve efficiency if trade is always efficient. A cooperative investment would for example include situations in which the seller's investment affects the buyer's consumption utility. Our model can be viewed as a general representation of this case.

institutions, such as firms. The exact mechanism by which this occurs is a question for future research.

## References

- Binmore, K., C. Proulx, and L. Samuelson (1995). Hard bargains and lost opportunities. Technical Report 9517, University of Wisconsin.
- Binmore, K. G., A. Rubinstein, and A. Wolinsky (1986). The nash bargaining solution in economic modeling. *Rand Journal of Economics* 17, 176–88.
- Carmichael, H. L. and W. B. MacLeod (1999). Caring about sunk costs: A behavioral solution to hold-up problems with small stakes. Olin Working Paper 99-19, University of Southern California.
- Che, Y.-K. and D. B. Hausch (1999, March). Cooperative investments and the value of contracting. *American Economic Review* 89(1), 125–47.
- Cole, H. L., G. J. Mailath, and A. Postlewaite (2000). Efficient non-contractible investments in a finite economy. Mimeo, University of Pennsylvania.
- Ellingsen, T. and J. Robles (2000). Does evolution solve the hold-up problem? Working Paper.
- Ellison, G. (2000). Basins of attraction, long run stability and the speed of step-by-step evolution. *Review of Economic Studies* 67, 17–45.
- Felli, L. and K. Roberts (2000). Does competition solve the hold-up up problem? Technical Report, London School of Economics.
- Freidlin, M. and A. Wentzell (1984). *Random Perturbations of Dynamical Systems*. Berlin: Springer.
- Grossman, S. J. and O. D. Hart (1986, August). The costs and benefits of ownership: A theory of vertical and lateral integration. *Journal of Political Economy* 94(4), 691–719.
- Grout, P. (1984, March). Investment and wages in the absence of binding contracts: A nash bargaining approach. *Econometrica* 52(2), 449–460.
- Hart, O. and J. Moore (1999, January). Foundations of incomplete contracts. *Review of Economic Studies* 66(1), 115–138.
- Hart, O. D. and J. H. Moore (1990). Property rights and the nature of the firm. *Journal of Political Economy* 98, 1119–58.
- Kandori, M., G. Mailath, and R. Rob (1993). Learning, mutation and long run equilibria. *Econometrica* 61, 27–56.
- MacLeod, W. B. and J. M. Malcomson (1993, September). Investments, holdup, and the form of market contracts. *American Economic Review* 83(4), 811–837.
- Robles, J. (2001). The evolution of contracts and property rights. Technical Report 11, University of Colorado.
- Thaler, R. (1980). Toward a positive theory of consumer choice. *Journal of Economic Behavior and Organization* 1, 39–60.

Tröger, T. (2000). Why sunk costs matter for bargaining outcomes: An evolutionary approach. mimeo, University College London.

Young, H. P. (1993). An evolutionary model of bargaining. *Journal of Economic Theory* 59, 145–168.

## Appendix

### Proof of Proposition 1:

Efficiency implies  $S_{HH} - c > S_{LH} > 0$ , therefore if one sets  $\hat{x}_{LH} = 0$ , then conditions 1 and 3 for a stable convention are strictly satisfied for  $\lambda = 0$ . The trade condition implies that  $S_{LH} > \delta(S_{LL} + S_{HH})/2 > \delta S_{HH}/2$  and therefore condition 2 is strictly satisfied. Given that the expressions in the definition of stability are continuous for small  $\lambda$ , the conditions for stability are satisfied for small  $\lambda$ .  $\square$

### Proof of Lemma 1:

Let  $s$  and  $s'$  be two arbitrary states in  $\mathcal{S}$ . We show that there is a positive multi-step transition probability from  $s$  to  $s'$  and a positive one-step transition probability from  $s'$  to  $s'$ . This then implies that the process is irreducible and aperiodic.

Assume that  $\sigma_t = s$ . With positive probability the bargaining strategy of all agents at time  $t$  is such that all agents carry out the project (some mutations of bargaining strategies might be needed) and leave the population. Hence, with positive probability in period  $t + 1$  the types of all agents in the population are determined anew and with a positive probability the resulting distribution of types matches exactly the one in  $s'$ . Every period there is positive probability that the distribution of types stays like that. If there are both high and low types in  $s'$  it is straight-forward to see that any set of observations needed to create empirical distribution functions which have positive weight in  $s'$  can be created by multiple mutations of bargaining behavior of the agents given the type distribution. In case there are only high or only low types in  $s'$  consider the transition where first all but one agent get the type required in  $s'$ , then all the observations needed to create all the beliefs in  $s'$  are created by mutations and finally the single agent with a different type leaves the population and changes her type. In any case there is a positive probability that  $s'$  is reached in multiple steps. Furthermore, since there is always a positive probability that all agents only observe matches between the same types during a period and therefore do not change their beliefs, there is a positive probability that the process stays in  $s'$  once it has reached  $s'$ . Hence, the process is irreducible and aperiodic.  $\square$

### Proof of Lemma 2:

**a:**

First we show that all the sets given in the Lemma are limit sets, i.e. we have to show that for  $\epsilon = 0$  they are absorbing and for each pair of states in such a set there is a positive (multi-step) transition probability.

It follows from the definition of  $\mathcal{C}(\delta)$  that if  $x \in \mathcal{C}(\delta)$  and all individuals have point beliefs  $\beta$  such that  $\hat{F}_{LH}(\cdot, \beta) = \mathcal{P}(x)$ ,  $\hat{F}_{HL}(\cdot, \beta) = \mathcal{P}(S_{LH} - x)$ , all individuals have the optimal bargaining strategy

$x_{LH} = x, x_{HL} = S_{LH} - x$ . Therefore, in the absence of mutations these point beliefs can never be altered and therefore  $\Omega(x)$  is absorbing. Furthermore, since in every period every distribution of types has a positive probability regardless of the actual investment behavior, and so also for every  $\hat{p} \in P$  there is a positive probability that a sample yielding such an estimator is observed, all possible distributions of types and  $\hat{p}$  can be reached with positive probability. Hence, the set  $\Omega(x)$  is connected, which implies that it is a limit set.

To proof that these are the only limit sets, we show that from every state which is not in one of the limit sets described above there is a positive probability to reach one of these sets. This comes down to showing that a homogeneous bargaining convention which is consistent with all  $\hat{p} \in P$  can always be reached with positive probability. The transition can go as follows: assume  $\sigma_t = s$  for some arbitrary state  $s \in \mathcal{S}$ . With positive probability there are at least  $m$  low types in  $\sigma_{t+1}$  and with positive probability at  $t + 2$  there is some pairing of a low type agent  $a_L$  and a high type agent  $a_H$  with bids  $\tilde{x}_{HL}, \tilde{x}_{LH}$ , where  $a_L$  has beliefs  $\beta$  such that  $\hat{p}(\beta) = 0$  and accordingly  $\tilde{x}_{LH} \geq \frac{\delta S_{LL}}{2}$ . With positive probability this pairing is repeated  $m$  times from period  $t + 2$  till  $t + m - 1$  and one agent, we call him  $b_H$ , in the population samples all these pairings but no other high-low pairings. Accordingly, at  $t + m$  she has beliefs such that  $\hat{F}_{LH}(\cdot, \beta_{t+m}) = \mathcal{P}(\tilde{x}_{LH})$ . Furthermore, there is a positive probability that the beliefs of  $a_L$  (or the agent who replaces her) only observes high-high meetings during this period and her beliefs stay unchanged. Furthermore there is a positive probability that  $a_L$  and  $b_H$  are matched in periods  $t + m$  till  $t + 2m - 1$ . In each such matching the two bids are  $\tilde{x}_{LH}$  of  $a_L$  and  $S_{LH} - \tilde{x}_{LH}$  of  $b_H$ . Again, there is a positive probability that all individuals sample only these high/low pairings during periods  $t + m$  to  $t + 2m - 1$ . Then in  $t + 2m$  all agents have beliefs such that  $\hat{F}_{LH}(\cdot, \beta_{t+2m}) = \mathcal{P}(\tilde{x}_{LH}), \hat{F}_{HL}(\cdot, \beta_{t+2m}) = \mathcal{P}(S_{LH} - \tilde{x}_{LH})$ . If  $S_{LH} - \tilde{x}_{LH} \geq \frac{\delta S_{HH}}{2}$  we have  $\tilde{x}_{LH} \in \mathcal{C}(\delta)$  and the proof of (a) is complete.

If  $S_{LH} - \tilde{x}_{LH} < \frac{\delta S_{HH}}{2}$ , there is a positive probability that in period  $t + 2m + 2$  there is a high type with  $\hat{p} = 1$ . This agent then makes a bid  $\tilde{x}_{HL}$  such that  $\tilde{x}_{HL} - \alpha < \frac{\delta S_{HH}}{2} \leq \tilde{x}_{HL}$  and the same arguments as above imply that there is a positive probability that a homogeneous state will evolve where for all agents hold beliefs  $\beta$  such that  $\hat{F}_{LH}(\cdot, \beta) = \mathcal{P}(S_{LH} - \tilde{x}_{HL}), \hat{F}_{HL}(\cdot, \beta) = \mathcal{P}(\tilde{x}_{HL})$ . Since  $\delta < \frac{2S_{LH}}{S_{HH} + S_{LL}}$  implies  $\frac{\delta S_{HH}}{2} > S_{LH} - \frac{\delta S_{LL}}{2}$ , we have  $S_{LH} - \tilde{x}_{HL} \in \mathcal{C}(\delta)$  for sufficiently small  $\alpha$ .

**b:**

Assume  $\sigma_t = s$  for an arbitrary state  $s \in \mathcal{S}$ . Assume that there are at least  $m$  low types and at least  $m$  high types in the population (if this is not true, there is a positive probability that at least  $m$  low and high types will be in the population within two periods). Then, there is a positive probability that in period  $t + 1$  all low types have beliefs  $\hat{p}^i = 0$  and at least  $m$  are matched with high types. The resulting demands at  $t + 1$  of these low types are larger or equal to  $\underline{x}_{LH}(0)$ . There is a positive probability that at least  $m$  high types observe these  $m$  demands in  $t + 2$  and that the same  $m$  high types in period  $t + 3$  have beliefs such that  $\hat{p}(\beta^i) = 1$  and are matched with low types. Since for these individuals we have  $\hat{F}_{LH}(\cdot, \beta) = \mathcal{P}(\underline{x}_{LH}(0))$  and  $\underline{x}_{LH}(0) > \bar{x}_{LH}(1)$  all the outside option is binding for all these high types and they demand  $x_{HL} = S_{LH} - \bar{x}_{LH}(1)$  in period  $t + 3$ . With positive probability these  $m$  demands are sampled by all agents in  $t + 4$  and hence all agents have beliefs such that  $\hat{F}_{HL}(\cdot, \beta) = \mathcal{P}(S_{LH} - \bar{x}_{LH}(1))$ . With positive probability these beliefs stay unchanged till  $t + 5$  whereas the belief about the type distribution changes to

$\hat{p}(\beta) = 0$ . With positive probability in  $t + 5$  now at least  $m$  low type agents are matched with high types and since  $\underline{x}_{LH}(0) > \bar{x}_{LH}(1)$  their outside option is binding and their demands are  $x_{LH} = \underline{x}_{LH}(0)$ . With positive probability all agents sample the demands of these  $m$  low types in  $t + 6$  and hence all agents have beliefs  $\beta$  such that  $(\hat{p}(\beta) = 0, \hat{F}_{HL}(\cdot, \beta) = \mathcal{P}(S_{LH} - \bar{x}_{LH}(1)), \hat{F}_{LH}(\cdot, \beta) = \mathcal{P}(\underline{x}_{LH}(0)))$ . We denote this state by  $\tilde{s}$ . The fact that there exists a positive multi-step transition probability from every state to  $\tilde{s}$  implies that the Markov chain has a single limit set which includes  $\tilde{s}$ . Obviously, this single limit set consists of all states which can be reached with positive probability from  $\tilde{s}$ . Taking into account that every demand of a high type where the outside option is binding has to be in  $\mathcal{B}_{HL}$  and that the best response of a high type with some beliefs  $\hat{F}_{LH}$  with support in  $\mathcal{B}_{LH}$  and  $\hat{p} \in P$  must lie in  $\mathcal{B}_{HL}$  as well, shows that all demands of high types have to be in  $\mathcal{B}_{HL}$  once  $\tilde{s}$  has been reached. Similarly for a low type. Accordingly, given that  $\epsilon = 0$ , any observation outside  $\mathcal{B}_{HL} \times \mathcal{B}_{LH}$  has probability zero once  $\tilde{s}$  has been reached before.  $\square$

### Proof of Proposition 2:

We have to determine which of the limit sets characterized in Lemma 2 are stochastically stable. We use the radius modified coradius criterion introduced in Ellison (2000). For a union of limit sets  $\Omega$  the radius  $R(\Omega)$  is defined as the minimum number of mutations needed to get to a state outside the basin of attraction of  $\Omega$  with positive probability. The modified coradius  $CR^*(\Omega)$  is defined as follows: consider an arbitrary state  $x \notin \Omega$  and a path  $(z_1, z_2, \dots, z_T)$  from  $x$  to  $\Omega$  where  $L_1, L_2, \dots, L_r \subset \Omega$  is the sequence of limit sets the path goes through (this implies  $L_r \subseteq \Omega$ ). We define the modified costs of this path by

$$c^*(z_1, \dots, z_T) = c(z_1, \dots, z_T) - \sum_{i=2}^{r-1} R(L_i),$$

where  $c(z_1, \dots, z_T)$  gives the number of mutations needed on the path  $(x_1, \dots, z_T)$ . Denoting by  $c^*(x, \Omega)$  the minimal modified costs for all paths from  $x$  to  $\Omega$  we define the modified coradius as

$$CR^*(\Omega) = \max_{x \notin \Omega} c^*(x, \Omega).$$

Ellison (2000) proves that every union of limit sets  $\Omega$  with  $R(\Omega) < CR^*(\Omega)$  contains all stochastically stable states.

In what follows we calculate the radius and modified coradius of the bargaining conventions described in Lemma 2. In the case of substitutes the limit sets are of the form  $\Omega(x_{LH})$  for  $x_{LH} \in \mathcal{C}(\delta)$ . Let  $\tilde{x}_{LH}$  be an arbitrary bargaining convention with  $\tilde{x}_{LH} \in \mathcal{C}(\delta)$ . To destabilize the convention upwards either a sufficient number of high types have to mutate to a  $x_{HL}$  smaller than  $S_{LH} - \tilde{x}_{LH}$ , in the extreme case  $x_{HL} = 0$ , such that the best response of a high type who has sampled all these mutants becomes  $x_{LH} = S_{LH}$ , or a sufficient number of low types have to mutate to  $x_{LH} = \tilde{x}_{LH} + \alpha$  such that the best response of a high type who has sampled all these mutants becomes  $x_{HL} = S_{LH} - \tilde{x}_{LH} - \alpha$ , where  $\alpha = \frac{S_{LH}}{k}$ . As has been demonstrated in Young (1993), for sufficiently small  $\alpha$  the second of these two possibilities yields transitions with a lower number of mutations (the number goes to zero as  $\alpha$  goes to zero). Similar arguments hold for a downwards destabilization and therefore in order to leave a convention  $\tilde{x}_{LH}$  with the minimal necessary number of mutations either the path to  $\tilde{x}_{LH} + \alpha$  or the path to  $\tilde{x}_{LH} - \alpha$  has to be taken. We define by

$c_+(x_{LH})$  the minimal number of mutations needed to get to  $\tilde{x}_{LH} + \alpha$  and by  $c_-(x_{LH})$  the minimal number of mutations needed to get to  $\tilde{x}_{LH} - \alpha$ . We first calculate  $c_+(\tilde{x}_{LH})$ .

The number of mutations needed to destabilize a convention also depends on the beliefs  $\hat{p}$ . We first show that the minimal number of mutants either occurs at  $\hat{p} = 0$  or at  $\hat{p} = 1$ . Consider a low type whose beliefs  $\hat{F}_{HL}$  attach probability  $q$  to  $x_{HL} = S_{LH} - \tilde{x}_{LH} + \alpha$  and  $1 - q$  to  $x_{HL} = S_{LH} - \tilde{x}_{LH}$ . Denote by  $v$  the expected discounted payoff of this individual given that he faces a high type and bids  $x_{LH} = \tilde{x}_{LH}$  whenever facing a high type. Taking into account that he will always trade immediately when he meets another low type we get

$$v = (1 - q)\tilde{x}_{LH} + \delta q \left( \hat{p}v + (1 - \hat{p})\frac{S_{LL}}{2} \right)$$

and

$$v(q; \hat{p}) := \frac{(1 - q)\tilde{x}_{LH} + \delta q(1 - \hat{p})S_{LL}/2}{1 - \delta q\hat{p}}.$$

Note that this expression is monotonous in  $\hat{p}$  for  $\hat{p} \in [0, 1]$  (increasing or decreasing). The minimal number of mutations needed to destabilize the convention is given by  $\lceil m\tilde{q} \rceil$ , where  $\tilde{q}$  is the minimal  $q$  such that:

$$v(q; \hat{p}) < \tilde{x}_{LH} - \alpha$$

holds for some  $\hat{p} \in [0, 1]$ . Since the right hand side is constant in  $q$  and  $\hat{p}$  and the left hand side is monotonous in  $\hat{p}$  for all  $q$  the minimal  $q$  is either attained at  $\hat{p} = 0$  or at  $\hat{p} = 1$ .

With  $\hat{p} = 0$  we get

$$v(q; 0) = (1 - q)\tilde{x}_{LH} + \delta q\frac{S_{LL}}{2},$$

which gives

$$q > q_{1-}(\tilde{x}_{LH}) := \frac{\alpha}{\tilde{x}_{LH} - \delta\frac{S_{LL}}{2}}.$$

For  $\hat{p} = 1$  we have

$$v(q, 1) = \frac{1 - q}{1 - \delta q}\tilde{x}_{LH}.$$

Accordingly, the convention can be destabilized downwards if

$$q < q_{2-}(\tilde{x}_{LH}) := \frac{\alpha}{\tilde{x}_{LH}(1 - \delta) + \delta\alpha}.$$

Comparing the two we see that  $q_{1-}(\tilde{x}_{LH}) < q_{2-}(\tilde{x}_{LH})$  if and only if  $\tilde{x}_{LH} > \frac{S_{LL}}{2} + \alpha$ . All-together we have

$$c_-(\tilde{x}_{LH}) = \begin{cases} q_{1-}(\tilde{x}_{LH}) & \tilde{x}_{LH} \geq \frac{S_{LL}}{2} + \alpha \\ q_{2-}(\tilde{x}_{LH}) & \tilde{x}_{LH} < \frac{S_{LL}}{2} + \alpha. \end{cases}$$

Similar reasoning for destabilizations upwards shows for a high type, who is matched with a low type and who believes that a fraction  $q$  of low types demands  $x_{LH} = \tilde{x}_{LH} + \alpha$  and a fraction  $1 - q$  of low types demands  $x_{LH} = \tilde{x}_{LH}$ , has the following expected payoff from demanding  $x_{HL} = S_{LH} - \tilde{x}_{LH}$ :

$$\begin{aligned} w(q; 0) &= \frac{1 - q}{1 - \delta q}(S_{LH} - \tilde{x}_{LH}) \\ w(q, 1) &= (1 - q)(S_{LH} - \tilde{x}_{LH}) + \delta q\frac{S_{HH}}{2}. \end{aligned}$$



This implies

$$c_+(\tilde{x}_{LH}) = \begin{cases} q_{1+}(\tilde{x}_{LH}) & \tilde{x}_{LH} \geq S_{LH} - \frac{S_{HH}}{2} - \alpha \\ q_{2+}(\tilde{x}_{LH}) & \tilde{x}_{LH} < S_{LH} - \frac{S_{HH}}{2} - \alpha, \end{cases}$$

where

$$q_{1+} = \frac{\alpha}{(S_{LH} - \tilde{x}_{LH})(1 - \delta) + \delta\alpha}$$

$$q_{2+} = \frac{\alpha}{S_{LH} - \tilde{x}_{LH} - \delta\frac{S_{HH}}{2}}.$$

The function  $c_-$  is decreasing in  $\tilde{x}_{LH}$  whereas  $c_+$  is increasing in this variable which implies that they have a unique intersection. We denote this intersection point by  $\hat{x}_{LH}$ . Clearly at this point  $\min[c_-, c_+]$  is maximized. For

$$(7) \quad \delta \leq \frac{2(S_{HH} - S_{LH})}{S_{HH} - S_{LL}}$$

$\hat{x}_{LH}$  lies on the intersection of  $q_{1-}$  and  $q_{1+}$  and is given by

$$(8) \quad \hat{x}_{LH} = \frac{S_{LH}}{2} - \frac{\delta}{2(2 - \delta)}(S_{LH} - S_{LL} - 2\alpha).$$

To establish (a) we first observe that under the assumptions made in (a) the condition (7) holds and  $\hat{x}_{LH} \in [\frac{\delta S_{LL}}{2}, S_{LH} - \frac{\delta S_{HH}}{2}]$  for small  $\alpha$ . Hence, there exists a  $\hat{\hat{x}}_{LH} \in \mathcal{C}(\delta)$  that maximizes  $\min[c_+, c_-]$  over  $\mathcal{C}(\delta)$  and whose distance from  $\hat{x}_{LH}$  is smaller than  $\alpha$ . Taking into account Lemma 2 this in particular implies that there is a limit set corresponding to the bargaining convention  $\hat{\hat{x}}_{LH}$ .

>From the arguments above it follows that for every  $x_{LH} \in \mathcal{C}(\delta)$  with  $x_{LH} < \hat{\hat{x}}_{LH}$  we have for the radius of the limit set  $\Omega(x_{LH})$ :  $R(\Omega(x_{LH})) = \lceil mc_+(x_{LH}) \rceil$  and for every  $x_{LH} \in \mathcal{C}(\delta)$  with  $x_{LH} > \hat{\hat{x}}_{LH}$  we have  $R(\Omega(x_{LH})) = \lceil mc_-(x_{LH}) \rceil$ . >From every limit set  $\Omega(x_{LH})$  there is a path to  $\Omega(\hat{\hat{x}}_{LH})$  along a graph  $g$  which connects every limit set  $\Omega(x_{LH})$  where  $x_{LH} < \hat{\hat{x}}_{LH}$  with  $\Omega(x_{LH} + \alpha)$ , and every limit set  $\Omega(x_{LH})$  where  $x_{LH} > \hat{\hat{x}}_{LH}$  with  $\Omega(x_{LH} - \alpha)$ . This implies that

$$CR^*(\Omega(\hat{\hat{x}}_{LH})) \leq \max_{x_{LH} \in \mathcal{C}(\delta) \setminus \{\hat{\hat{x}}_{LH}\}} R(\Omega(x_{LH})).$$

For sufficiently large  $m$  we have  $R(\Omega(\hat{\hat{x}}_{LH})) > R(\Omega(x_{LH}))$  for all  $x_{LH} \in \mathcal{C}(\delta) \setminus \{\hat{\hat{x}}_{LH}\}$  and therefore  $R(\Omega(\hat{\hat{x}}_{LH})) > CR^*(\Omega(\hat{\hat{x}}_{LH}))$ . Using the radius-modified coradius criterion we can conclude that the limit set corresponding to  $\hat{\hat{x}}_{LH}$  is stochastically stable. For  $k \rightarrow \infty$  we have  $\hat{\hat{x}}_{LH} \rightarrow \hat{x}_{LH}$  and get (a). Exactly the same arguments establish (b), where it has to be taken into account that in this case  $\hat{x}_{LH}$  lies at the intersection of  $q_{1-}$  and  $q_{2+}$  which is given by

$$\hat{x}_{LH} = \frac{S_{LH}}{2} - \frac{\delta}{4}(S_{HH} - S_{LL})$$

□

**Proof of Lemma 3:**

Parts (a) and (b) are trivial. To prove part (c) we denote by  $Q(\lambda) = [q_{ij}(\lambda)]_{i,j \in \tilde{S}}$  the one-step transition matrix of the Markov process  $\{\tilde{\sigma}_t\}$ . It can then easily be established that  $\lim_{\lambda \rightarrow 0} q_{i0} + q_{i1} > 0$  for all  $i \in \tilde{S}$ . Furthermore, at state  $i = 0$  no individual can sample any high types and hence we have  $\hat{p}(\beta) = 0$  for all individuals and accordingly all choose high investment. Therefore  $\lim_{\lambda \rightarrow 0} q_{01} = 1$  and by the same reasoning  $\lim_{\lambda \rightarrow 0} q_{10} = 1$ . Therefore, the only limit set for  $\lambda \rightarrow 0$  is  $\{0, 1\}$  which implies that  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_i^*(\lambda) = 0$  for all  $i \in \tilde{S} \setminus \{0, 1\}$ . Using this we get from the Chapman-Kolmogoroff equation at state 0

$$\tilde{\pi}_0^*(\lambda) \left( \sum_{i \in \tilde{S} \setminus \{0\}} q_{0,i} \right) = \sum_{i \in \tilde{S} \setminus \{0\}} q_{i,0} \tilde{\pi}_i^*(\lambda)$$

that  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_0^* = \lim_{\lambda \rightarrow 0} \tilde{\pi}_1^* = 0.5$ . □

#### Proof of Proposition 5:

The proof of (b) is identical to the proof of part (a) of proposition 2. To proof (a) we again follow the proof of proposition 2 but observe that for  $S_{LH} < \frac{\delta}{2}((2 - \delta)S_{HH} + \delta S_{LL})$  we have  $\hat{x}_{LH} > S_{LH} - \delta \frac{S_{HH}}{2}$ . Therefore the point which maximizes  $\min[c_+, c_-]$  over  $\mathcal{C}(\delta)$  is given by  $\hat{x}_{LH}$  where  $\hat{x}_{LH} \leq S_{LH} - \delta \frac{S_{HH}}{2} < \hat{x}_{LH} + \alpha$ . Stochastic stability of the limit set  $\Omega(\hat{x}_{LH})$  is established analogous to the proof of proposition 2 but here we have  $\hat{x}_{LH} \rightarrow S_{LH} - \delta \frac{S_{HH}}{2}$  for  $k \rightarrow \infty$ . □

#### Proof of Lemma 4:

We show the proposition for  $p^*(x_{LH}) > 0.5$ , the other case analogous.

We denote again the one-step transition matrix of the process  $\{\tilde{\sigma}_t\}$  by  $Q = [q_{ij}(\lambda)]_{i,j \in \tilde{S}}$ . We can write these transition probabilities as

$$q_{ij} = \binom{n}{j} \beta_i^j (1 - \beta_i)^{n-j},$$

where

$$\beta_i = (1 - \lambda)s(mp^*(x_{LH}); i) + \lambda(1 - s(mp^*(x_{LH}); i))$$

is the probability that a randomly chosen individual is of high type. Note that for a given bargaining convention the investment decision only depends on the number of high types sampled by an individual in the current period. We denote by

$$s(mp^*(x_{LH}); i) = \sum_{k \geq mp^*} \binom{m}{k} \left(\frac{i}{n}\right)^k \left(1 - \frac{i}{n}\right)^{m-k}$$

the probability that an individual samples more than  $mp^*(x_{LH})$  high types in a population with  $i$  high types. Since we are dealing with the case of investment complements here, this is the probability of high investment.

Note first that

$$\begin{aligned}
1 - s(mp^*; i) &= \sum_{k < mp^*} \binom{m}{k} \left(\frac{i}{n}\right)^k \left(1 - \frac{i}{n}\right)^{m-k} \\
&= \sum_{k > m(1-p^*)} \binom{m}{k} \left(\frac{i}{n}\right)^{m-k} \left(1 - \frac{i}{n}\right)^k \\
&> \sum_{k > mp^*} \binom{m}{k} \left(1 - \frac{n-i}{n}\right)^{m-k} \left(\frac{n-i}{n}\right)^k \\
&= s(mp^*; n-i),
\end{aligned}$$

where the inequality follows from  $p^* > 0.5$ . Using this we get that for  $\lambda < 0.5$

$$\begin{aligned}
1 - \beta_{n-i} &= 1 - (1-\lambda)s(mp^*; n-i) - \lambda(1 - s(mp^*; n-i)) \\
&= 1 - \lambda - s(mp^*; n-i)(1-2\lambda) \\
&> 1 - \lambda - (1 - s(mp^*; i))(1-2\lambda) \\
&= (1-\lambda)s(mp^*; i) + \lambda(1 - s(mp^*; i)) \\
&= \beta_i.
\end{aligned}$$

This means that in a population with  $i$  high types the probability that an individual becomes a high type is smaller than the probability that an individual becomes a low type in a population with  $i$  low types. In particular, this implies that the probability that at least  $z$  individuals become high types in state  $\frac{i}{n}$  is smaller than the probability that at least  $z$  individuals become low types in state  $\frac{n-i}{n}$  for all  $z$ . We denote by  $\mathcal{L} = \{0, \frac{1}{n}, \dots, \frac{n-2}{2n}\}$ ,  $\mathcal{H} = \{\frac{n+2}{2n}, \dots, \frac{n-1}{n}, 1\}$ ,  $\tilde{\mathcal{L}} = \mathcal{L} \cup \{\frac{n}{2}\}$  and by  $\tilde{\mathcal{H}} = \mathcal{H} \cup \{\frac{n}{2}\}$ . Furthermore we denote by  $q_{iL}$  the transition probability from state  $i$  into the set  $\mathcal{L}$  and analogous the transition probabilities into the other sets defined above. The arguments above imply that

$$q_{iL} > q_{n-iH} \quad \text{and} \quad q_{i\tilde{L}} > q_{n-i\tilde{H}} \quad \forall i.$$

Note that both  $(\mathcal{L}, \tilde{\mathcal{H}})$  and  $(\tilde{\mathcal{L}}, \mathcal{H})$  are partitions of the state space, therefore under the stationary distribution the flows between  $\mathcal{L}$  and  $\tilde{\mathcal{H}}$  must be identical in both directions and so have to be the flows between  $\tilde{\mathcal{L}}$  and  $\mathcal{H}$ . This gives

$$(9) \quad \sum_{i=0}^{n/2-1} \tilde{\pi}_i^* q_{i\tilde{H}} = \sum_{i=0}^{n/2} \tilde{\pi}_{n-i}^* q_{n-iL}$$

$$(10) \quad \sum_{i=0}^{n/2} \tilde{\pi}_i^* q_{iH} = \sum_{i=0}^{n/2-1} \tilde{\pi}_{n-i}^* q_{n-i\tilde{L}}$$

Our claim can now be easily shown by contradiction. If  $\sum_{i=0}^{n/2-1} \tilde{\pi}_i^* \leq \sum_{i=0}^{n/2-1} \tilde{\pi}_{n-i}^*$  we can use  $q_{iH} < q_{n-iL}$  for all  $i$  to derive that

$$\sum_{i=0}^{n/2} \tilde{\pi}_i^* q_{iH} < \sum_{i=0}^{n/2} \tilde{\pi}_{n-i}^* q_{n-iL}$$

Since (9, 10) have to hold this implies

$$\sum_{i=0}^{n/2-1} \tilde{\pi}_{n-i}^* q_{n-i\tilde{L}} < \sum_{i=0}^{n/2-1} \tilde{\pi}_i^* q_{i\tilde{H}}$$

But we also have  $q_{i\tilde{H}} < q_{n-i\tilde{L}}$  for all  $i$  and therefore this inequality contradicts our assumption that  $\sum_{i=0}^{n/2-1} \tilde{\pi}_i^* \leq \sum_{i=0}^{n/2-1} \tilde{\pi}_{n-i}^*$ . Accordingly, we must have  $\sum_{i=0}^{n/2-1} \tilde{\pi}_i^* > \sum_{i=0}^{n/2-1} \tilde{\pi}_{n-i}^*$

To show that  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_0^* = 1$  we can again apply the radius-modified coradius criterion. For  $\lambda = 0$  there are two limit sets, namely  $\{0\}$  and  $\{1\}$ . In order to invest high an individual has to sample at least  $\lceil mp^* \rceil$  high types. Therefore the radius of  $\{0\}$  is given by  $R(\{0\}) = \lceil mp^* \rceil$ . On the other hand, the state where maximal the number of mutations is needed to have a positive transition probability into  $\{0\}$  is the state 1 and therefore we have  $CR^*(\{0\}) = \lceil m - mp^* \rceil$ . For  $p^* > 0.5$  this implies that  $R(\{0\}) > CR^*(\{0\})$  for sufficiently large  $m$  and therefore  $\lim_{\lambda \rightarrow 0} \tilde{\pi}_0^* = 1$ .  $\square$