# Contracting Under Reciprocal Altruism

Shchetinin, Oleg

Toulouse School of Economics, Toulouse School of Economics - GREMAQ

16 February 2009

# Contracting under Reciprocal Altruism *

Oleg Shchetinin

Toulouse School of Economics (GREMAQ)

May 25, 2009

### Abstract

I develop a model of contracting under reciprocal altruism accounting for some evidence which is paradoxical from the point of view of neo-classical models with selfish actors. My model predicts the crowding-out effect observed in the Trust Game with the possibility of a fine; for the Control Game the model predicts that an equilibrium can exhibit "no effect of control", "hidden cost of control", or "positive effect of control", depending on the characteristics of the actors, as observed in the lab. This suggests that reciprocal altruism modeling could be fruitful more generally in applications of contract theory.

**Keywords:** Contract Theory, Signaling, Behavioral Economics, Reciprocal Altruism, Extrinsic and intrinsic motivation, Experimental Economics.

**JEL Classification Numbers:** D82, M54

## 1   Introduction

Intriguing observations about human response to incentives have recently been made. For instance, providing additional incentives can, in contrast with standard models with selfish actors, lead to lower levels of performance and intentions seem to matter, according to Fehr and Rockenbach [2003], Falk and Kosfeld [2006] and many others[1].

In this paper I develop a Principal-Agent model embodying reciprocal altruism. An important contribution of the paper is that a simple formal model of reciprocal altruism is able to give reliable predictions for human behavior. While the idea that reciprocity, altruism and other forms of social preferences shape people's behavior is not new, there are only a few models of reciprocal altruism in the literature.

My model is based on a simple idea that a person cares more about those who care more about him. In other words, a person is more altruistic towards

---

[1]See below the more detailed discussion of the relevant literature.

those whom he perceives as being altruistic towards him. This is the essence of the reciprocal altruism. In the Principal-Agent relationship, an altruistic Agent is inspired to exert effort even in the absence of monetary incentives, i.e. the Agent's altruism works as an intrinsic motivator. If furthermore, the Agent is reciprocal, the Principal should demonstrate his altruism in order to inspire the Agent's intrinsic motivation. This leads to a signaling game in which the Principal signals his altruism through offering a "generous" contract.

After analyzing the signaling game with intrinsic motivation only, I incorporate extrinsic incentives. I follow the frameworks of the two well known lab experiments of Fehr and Rockenbach [2003] and Falk and Kosfeld [2006] to investigate how the predictions of my model match the empirical evidence. I show that most of the observed patterns of behavior in the experiments are predicted by the model.

The behavioral patterns observed in the experiments are not limited to the lab, they are common in human relations, in the workplace in particular. The evidence from the field is documented in Gneezy [2002], Falk [2007], Bolton and Ockenfels [2008], Paarsch and Shearer [2007], Shearer [2004], Bellemare and Shearer [2007], Berry and Kanouse [1987], Maréchal and Thöni [2007].[2]

The relevance of altruism and reciprocity are supported by compelling evidence in the literature.

The Dictator Game introduced in Kahneman et al. [1986] provides evidence for pure altruism[3]. The survey of Andreoni [2006] demonstrates that impure altruism - taste for "warm glow" shapes people's decisions in many circumstances. The evidence for reciprocity[4] is provided by Berg et al. [1995] who introduced the Trust Game, which was repeated with modifications in Fehr and Rockenbach [2003], Fehr and List [2004] and others. The evidence also comes from variants of the Gift Exchange Game by Fehr et al. [1997]; Lost Wallet Game by Dufwenberg and Gneezy [2000] and Charness et al. [2007]; Moonlighting Game by Abbink et al. [2000].

Some recent papers develop models which can explain puzzling behavior through taking into account the interaction between extrinsic incentives[5] and intrinsic motivation[6], which can lead, for instance, to motivation crowding-out.

---

[2]However, Kube et al. [2006] found support for negative reciprocity and question positive reciprocity, especially in the long-run. Gneezy and List [2006] found reciprocity in short-run (the first 2 hours of work) and decreasing reciprocity in the long-run: to the end of the 6-hours job the subjects with more generous wage didn't work harder that the others. Some studies question the relevance of the lab experiments - see, e.g. List [2007], Hennig-Schmidt et al. [2005], List and Levitt [2005]. We should be warned by these studies but evidence for reciprocity comes from many different sources, so it's hard to question that reciprocity is an important psychological characteristic of human beings.

[3]In the various Dictator Game experiments, the subjects are endowed with a sum of money. They decide then on how much of this windfall endowment to give to a stranger. More than half of the subjects give between 20% and 50% of the endowment.

[4]Precisely, "intrinsic reciprocity", not "consequentialism" or "strategic reciprocity".

[5]The list of extrinsic motivators is not limited to the incentive payments (piece-rate wage or bonus payment) but includes also expectation of future material payoff e.g. reputation building due to long-term interaction, strategic reciprocity, career concerns, comparative performance based payment (tournaments), monitoring/control etc.

[6]The literature provides evidence for many kinds of intrinsic motivation, apart from altruism and reciprocity. The Ultimatum Game introduced by Güth et al. [1982] illustrates that taste for fairness and/or inequality aversion is an important factor determining behavior; another evidence for fairness comes from different versions of the Gift Exchange Game - see Fehr et al. [1993], Fehr and Falk [1999]. Social norms (avoiding social disapproval/geting

My paper does the similar, though it is based on a different mechanism. I show that reciprocal altruism itself can account for many systematically observed behavior patterns given a natural information structure.

Levine [1998] developed the model in which Agent's altruism is conditional on the Principal's one and applied it to explain behavior observed in a number of lab experiments, for instance the Ultimatum Game and the Public Good contribution game.

Falk and Fischbacher [2006] develop a theory of reciprocity based on psychological games, i.e. with utilities of the actors depending not only on their material payoffs, but also on the perceived intentions of another player, i.e. on the 1-st and 2-nd order beliefs, following Rabin [1993]. The players' concern about the equitable outcome plays an important role in the analysis, in contrast with my model, which is based on reciprocal altruism. The model of Falk and Fischbacher [2006] can explain behavior in the Ultimatum Game, the Gift Exchange Game and some other experiments.

Ellingsen and Johannesson [2008] build a model which can account for the crowding-out effect. It is similar to mine in that they consider altruistic actors. One difference is that they also incorporate the taste for the social esteem (pride) in the utility function, i.e. second order belief. Another difference is that the actors in their model are unconditionally altruistic.

Sliwka [2007] develops a model which can also account for the crowding-out effect through a mechanism based on social norms, which is different from mine. There are reliable and unreliable Agents in the model. The reliable Agents follow the contract whereas the unreliable ones can deviate even if a contract is signed. As a consequence, the incentive scheme has to be high-powered for the unreliable Agents to perform at a high level. This leads to the fact that observing the high-powered incentive scheme, the Agent can learn that there is a social norm to be unreliable which can crowd-out his intrinsic motivation based on inherent reliability.

In a recent paper Dur [2008] analyzes a model based on reciprocal altruism, applied to the workplace relation and shows the importance of attention payed to the Agent by the Principal. The model follows Levine [1998] approach.

My paper proceeds as follows. Section 2 describes the Core Model and presents its general analysis, leading to the benchmark results. Section 3 analyzes in detail signaling under reciprocal altruism without any extrinsic incentives, i.e. the Core Model. Section 4 introduces the extrinsic incentives in the Core Model and deals with two altered versions of it, which match the settings of the well known lab experiments - the Trust Game with a possibility of a fine by Fehr and Rockenbach [2003] and the "Control Game" by Falk and Kosfeld [2006]. I show that there is a unifying framework for the Core Model and its two extensions with the same mechanism of motivation crowding out. Section 5 concludes.

---

social approval) influence economic decisions. People can change their behavior under peer pressure or have a taste for the social embeddedness. The evidence are provided by a variant of the Gift Exchange Game in Gächter and Falk [2002] and Third Party Punishment Game by Fehr and Fischbacher [2004]. A person may have taste for the others' belief about his motivation (or type) - see Rabin [1993], Falk and Fischbacher [2006] and Bénabou and Tirole [2006a]. The list of intrinsic motivators can be continued with self-learning, working on interesting/challenging task (in this case effort may not be costly (painfull), the job rather gives fun and higher effort increases utility) etc.

## 2　The Core Model

Consider the Principal-Agent relation with one employer - the Principal and one worker - the Agent. Assume that the Principal is altruistic towards the Agent and the Agent reciprocates the employer's altruism: if the Agent perceives the Principal to care about him, he becomes altruistic towards the Principal. The Principal offers a contract to the Agent.

Assume that output is equal to effort so that there is no moral hazard. The output is observable and verifiable and can be contracted upon.

Producing output is costly with cost function $C(q)$ satisfying standard assumptions - convexity and zero cost at zero output:

$$C'(q) > 0, \ C''(q) > 0 \text{ for } q > 0$$
$$C(0) = 0, \ C'(0) = 0$$

Let $B$ be the Agent's benefit from interacting with the Principal. The benefit can be psychological or a monetary payment from a third party[7].

For now, assume that the Agent doesn't respond to monetary incentives, beyond some subsistence level, that we normalize to zero. In the Core model analysis, we focus on the intrinsic motivation only. In the subsequent section, I introduce extrinsic incentives. The selfish utilities of the Principal and the Agent are then given by

$$v = q$$
$$u = B - C(q)$$

Let $\alpha$ be the degree of the Principal's altruism, $\widehat{\alpha}$ - the Agent's perception of the Principal's altruism, and $\beta$ be the intensity of the Agent's reciprocity (more generally, it can be treated as intensity of intrinsic motivation of any nature emerging from perceiving the Principal as "generous"). I assume that $0 \leq \alpha, \beta \leq 1$. The interaction term $\widehat{\alpha}\beta$ represents the Agent's altruism emerging as a result of reciprocating altruism of the Principal.

The utilities of the Principal and the Agent when the Agent produces output $q$ are given by

$$V(q, \alpha) = v + \alpha u = q + \alpha(B - C(q)) \tag{1}$$
$$U(q, \widehat{\alpha}, \beta) = u + \widehat{\alpha}\beta v = B - C(q) + \widehat{\alpha}\beta q \tag{2}$$

The contract can be a command - "produce $q$" or can give the Agent some flexibility - say, "produce any quantity $q \in [q_1, q_2]$".

Notice the difference with the standard Principal-Agent setup. The Principal's valuation of the output is not always increasing, now it has an inverted-U shape: it increases for small enough values of output only and is maximal at $q = q^P$. Similarly, the Agent's payoff is not always decreasing in output, and has an inverted-U shape: it decreases only for large enough values and reaches the maximal value at $q^A$.

In what follows, I will refer to $q^P$ and $q^A$ as the Principal's and the Agent's preferred values of output (or performance). In contrast with the standard Principal-Agent models, $q^P \neq +\infty$, $q^A \neq 0$. Principal's and Agent's payoffs as functions of output are depicted in Figure 1.

---

[7]The latter is the case in the lab experiments which I consider in the paper. The third party will be an experimenter.
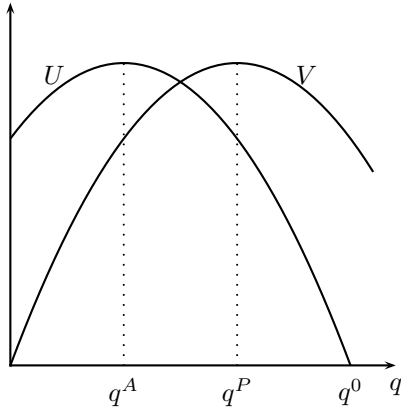
Figure 1: Principal's and Agent's payoffs under reciprocal altruism.

For $\alpha = \beta = 1$ the Principal's and the Agent's interests are aligned, $U(q) \equiv V(q)$, because there is full internalization, so that the two curves representing the Principal's and the Agent's utilities in Figure 1 coincide.

For smaller values of $\alpha$ or $\beta$, i.e. weaker internalization, there is a conflict of interest like in the standard Principal-Agent setup but this conflict is softened by the partial internalization of utilities. In the graph, the two inverted-U curves become more distant, and consequently, the distance between the maximizers of the Principal and Agent utilities $q^P$ and $q^A$ becomes larger: the Principal wants the Agent to exert more effort whereas the Agent prefers performing less.

For $\alpha$ and $\beta$ close to 1 the Agent's participation constraint is not binding because $q^P$ is "close enough" to the maximizer of the Agent's utility $q^A$ where the Agent utility is positive. However, as $\alpha$ or $\beta$ decrease, the participation constraint becomes binding.

Denote the value of $q$ starting from which the participation constraint is binding by $q^0(\alpha, \beta)$. I will refer to this value as the Agent's participation threshold.

## 2.1   Benchmark cases

The preferred value of output for the Principal is given by

$$q^P(\alpha) = \arg\max_q \left[ V(q, \alpha) \right] = \arg\max_q \left[ q - \alpha C(q) \right]$$

$$C'(q^P) = \frac{1}{\alpha} \tag{3}$$

If there are no barriers to implement this output level, such as Agent's participation constraint or limits on contract design, the Principal will induce it.

**Lemma 1.** *The Principal's preferred output $q^P(\alpha)$ is determined by (3) and is a decreasing function of $\alpha$:* $\frac{\partial q^P}{\partial \alpha} < 0$.

The lemma follows directly from (3).

5

The preferred value of output for the Agent is given by

$$q^A(\widehat{\alpha}, \beta) = \arg\max_q \left[ U(q; \widehat{\alpha}, \beta) \right] = \arg\max_q \left[ \widehat{\alpha}\beta q - C(q) \right]$$

$$C'(q^A) = \widehat{\alpha}\beta \tag{4}$$

This is output which is realized if the Agent is given full flexibility or if this level is available to the Agent despite some restrictions such as binding contract are imposed.

This equation shows that the Agent is willing to perform at an effort level such that his marginal cost is equal to his marginal benefit $\widehat{\alpha}\beta$. This justifies that $\widehat{\alpha}\beta$ can be considered as a measure of the Agent's intrinsic motivation.

**Lemma 2.** *The Agent's preferred output $q^A(\widehat{\alpha}, \beta)$ is determined by (4) and $q^A$ increases with $\alpha$ and $\beta$: $\frac{\partial q^A}{\partial \widehat{\alpha}} > 0$, $\frac{\partial q^A}{\partial \beta} > 0$.*

The lemma follows directly from (4).

For the case $\alpha, \beta < 1$ it's easy to see from (3) that $C'(q^P) > 1$, whereas (4) leads to $C'(q^A) < 1$, so that $q^P > q^A$ and there is always a gap between the Principal's and the Agent's preferred output levels which is larger for smaller $\alpha$ and $\beta$.

This leads to the following lemma.

**Lemma 3.** *The Principal's preferred output is always larger than the Agent's preferred output except for the case of $\alpha = \beta = 1$ when the preferred outputs are the same:*

$$q^P(\alpha) > q^A(\alpha, \beta) \quad \text{if } \alpha < 1 \text{ or } \beta < 1$$

$$q^P(\alpha = 1) = q^A(\alpha = 1, \beta = 1)$$

*For the case of symmetric or revealed information, so that $\widehat{\alpha} = \alpha$, the gap $q^P - q^A$ between the Principal's and the Agent's preferred outputs is a decreasing function of $\alpha, \beta$: $\frac{\partial \left( q^P - q^A \right)}{\partial \alpha} < 0$, $\frac{\partial \left( q^P - q^A \right)}{\partial \beta} < 0$*

The Agent's participation threshold $q^0$ is the unique root of the equation

$$U(q; \widehat{\alpha}, \beta) = B + \widehat{\alpha}\beta q - C(q) = 0 \tag{5}$$

The root exists and is unique since $U(0) = B > 0$, $U(q)$ increases for $q \in (0, q^A)$, so that $U(q^A) > 0$, then decreases for $q \in (q^A, \infty)$ and $U(q) \to -\infty$ as $q \to \infty$. Because of continuity of $U(q)$, there exists a unique $q^0 \in (q^A, \infty)$ such that $U(q^0) = 0$.

**Lemma 4.** *The Agent's participation threshold is given by $q^0(\widehat{\alpha}, \beta)$ which is an increasing function of $\widehat{\alpha}$ and $\beta$: $\frac{\partial q^0}{\partial \widehat{\alpha}} > 0$, $\frac{\partial q^0}{\partial \beta} > 0$*

The lemma is easy to obtain from (5).

I proceed now with the asymmetric information case analysis.

6

# 3  Signaling under Reciprocal Altruism

Consider the case when the Principal's altruism parameter $\alpha$ is her private information and the rest is symmetrically known.

The timing is as follows:

1. Principal learns $\alpha$.

2. Principal offers a contract[8]: a command, i.e. specifies the output $q$.

3. Agent accepts or rejects the contract.

4. Contract is implemented and payoffs are realized.

So, we have a signaling game with common values.

To make exposition even simpler, I consider the 2-type case and then generalize the result to the continuum-type case.

Even though the formal setting of the signaling game should be clear for most of the readers, I provide its formal description in the next subsection, in which some notation is introduced. I proceed then to the analysis of Perfect Bayesian Equilibria and refinement.

Readers who are not interested in the technical details can skip the technical subsection 3.1 and jump to 3.2 where the outcome of the signaling game as predicted by the refined equilibrium is described.

## 3.1  Signaling with 2 types

There are 2 players - Principal (sender) and Agent (receiver).

The Principal's type is her private information. Denote by $\mathcal{A}$ the set of the possible types, $\mathcal{A} = \{\alpha_H, \alpha_L\}$. The prior distribution is given by

$$\alpha = \begin{cases} \alpha_H & \text{with prob. } \Pi \\ \alpha_L & \text{with prob. } 1 - \Pi \end{cases} \tag{6}$$

The set of actions for the Principal[9] is $Q = [0, +\infty)$. The set of actions for the Agent is $A = [0, 1]$ with $a \in A$ be the probability of acceptance of an offer made by the Principal.

A pure strategy of the Principal is a type-contingent $q \in Q$, i.e. $q_H$ for $\alpha_H$ type and $q_L$ for $\alpha_L$-type.

A pure strategy of the Agent is an acceptance rule $a(\cdot)$. The value $a(q) \in [0, 1]$ is the probability of accepting the offer $q$. The set of Agent's pure strategies $\mathcal{F}$ is the set of all mappings from $Q$ to $[0, 1]$.

A mixed strategy of the Principal is a probability distribution over $Q$ conditional on type, $\sigma(\cdot | \alpha)$. Clearly her mixed strategy can be represented by the two probability distributions over $Q$: $(\sigma_H(\cdot), \sigma_L(\cdot))$.

A mixed strategy of the Agent is a probability distribution[10] $\tau$ over $\mathcal{F}$. The resulting mapping is still a mapping from $Q$ to $[0, 1]$. So, the mixed strategy of the Agent, denoted by $a_\tau(\cdot)$, is still an element of $\mathcal{F}$. I will will restrict attention to the pure strategies of the Agent.

---

[8]The contract is the take-it-or-leave-it offer.

[9]Actually, this set of actions can be reduced to $[0, q_H^0]$.

[10]There is an issue of measurability over the space of functions.

The Agent's ex-post belief on the Principal's type distribution $\mu(\cdot)$ is

$$\mu(q) = Prob(\alpha = \alpha_H | q)$$

There is one-to-one correspondence between Agent's belief $\mu(q)$ and Agent's ex-post expectation on the Principal's altruism

$$\widehat{\alpha}(q) = \mu(q)\alpha_H + (1 - \mu(q))\alpha_L \tag{7}$$

which is paralleled in the Agent's ex-post expected payoff[11]

$$U_\mu(q, a; \beta) = U(q, a, \widehat{\alpha}; \beta)$$

The pure strategies-beliefs profile is thus $((q_H, q_L), a(\cdot), \mu(\cdot))$. The mixed strategies-belief profile is $((\sigma_H(\cdot), \sigma_L(\cdot)), a(\cdot), \mu(\cdot))$.

The payoffs in the game for the pure strategies profile $((q_L, q_H), a(\cdot))$ are

$$V(q, a(\cdot), \alpha) = V(q, \alpha)a(q) = (q - \alpha C(q) + \alpha B)a(q)$$
$$U(q, a(\cdot), \alpha; \beta) = U(q, \alpha; \beta)a(q) = (B - C(q) + \alpha\beta q)a(q)$$

which parallel (1) and (2).

The payoffs for the mixed strategies can be determined in a standard way.

Denote by $q_H^P$ and $q_L^P$ the preferred output of the high and low altruistic Principals correspondingly. Formally, $q_H^P = q^P(\alpha_H)$, $q_L^P = q^P(\alpha_L)$. From lemma 1 we have

$$q_H^P < q_L^P \tag{8}$$

Intuitively, the Principal who cares more about the Agent wants him to work less. Intuitively, since marginal cost of effort increases and the Principal partially internalizes this cost, the one with stronger internalization prefers to have lower marginal cost $C'(q)$ because marginal benefit from output is constant (equal to 1).

Denote by $q_H^0$ and $q_L^0$ the Agent participation thresholds when he learns that the Principal type is $\alpha_H$ and $\alpha_L$ correspondingly. Formally, $q_H^0 = q^0(\alpha_H, \beta)$, $q_L^0 = q^0(\alpha_L, \beta)$.

Denote the participation threshold when there is no update on the Principal's type by $q_E^0$:

$$q_E^0 = q^0(E\alpha, \beta)$$

where

$$E\alpha = \Pi\alpha_H + (1 - \Pi)\alpha_L$$

According to Lemma 4,

$$q_L^0 < q_E^0 < q_H^0. \tag{9}$$

Intuitively, if the Principal's type is revealed to the Agent, then he is willing to exert more effort for those Principal who cares more about him. This is natural since the worker internalizes the benefits from output and intensity of the internalization is higher for the worker connected with more altruistic Principal.

---

[11]It is easy to see that $U_\mu(q, a(\cdot); \beta) = \mu(q)(B - C(q) + \alpha_H\beta q))a(q) + (1 - \mu(q))(B - C(q) + \alpha_L\beta q))a(q) = (B + \widehat{\alpha}\beta q - C(q))a(q)$.

### 3.1.1 The Perfect Bayesian Equilibrium

As is standard in the signaling games, the set of the Perfect Bayesian Equilibria (PBE) is large. In this part of the paper I characterize the structure of any PBE and proceed to refinement in section 3.1.2.

Consider a PBE of the signaling game. I will denote the equilibrium (pure) offers of the Principal of type $\alpha_H$ ($\alpha_L$) by $q_H^*$ ($q_L^*$), the equilibrium acceptance rule of the Agent by $a^*(\cdot)$. The belief supporting the equilibrium is $\mu^*(\cdot)$. So, the pure strategies-belief equilibrium profile is $((q_H^*, q_L^*), a^*(\cdot); \mu^*(\cdot))$. Similarly, the mixed strategies-belief equilibrium profile is $((\sigma_H^*(\cdot), \sigma_L^*(\cdot)), a^*(\cdot); \mu^*(\cdot))$.

Denote by $BR_\mu(q)$ the best response acceptance rule for the Agent with ex-post belief $\mu(\cdot)$:

$$BR_\mu(q) = \arg\max_{a \in [0,1]} U_\mu(q, a; \beta) \tag{10}$$

For a given offer $q$ the Agent calculates the value $\mu(q)$ according to his belief and solves the maximization program (10), which picks the value $a(q)$ for the acceptance rule $a(\cdot)$ at point $q$ only. The Best Response doesn't impose any restriction on values $a(\widetilde{q})$ for $\widetilde{q} \neq q$.

**Lemma 5.** *For any belief $\mu(\cdot)$, the Best Response acceptance rule is a threshold with the threshold value*

$$\widehat{q}(q) = q^0(\widehat{\alpha}(q), \beta) \tag{11}$$

$$BR_\mu(q) = \begin{cases} 1 & \text{if } q < \widehat{q}(q) \\ 0 & \text{if } q > \widehat{q}(q) \\ \text{any } a \in [0,1] & \text{if } q = \widehat{q}(q) \end{cases}$$

*where $\widehat{\alpha}(q)$ is given by (7)*

*For any belief $\mu(\cdot)$ and any offer $q$*

$$q_L^0 \leq \widehat{q}(q) \leq q_H^0 \tag{12}$$

*so that for any Best Response acceptance rule, the Agent accepts at least offers $q < q_L^0$ and rejects any offer $q > q_H^0$.*

**Corollary 1.** *The equilibrium acceptance rule $a^*(\cdot)$ is a threshold with threshold value given by (11)*

Proof is given in the Appendix.

Now I proceed to the analysis of the Principal equilibrium offer.

First, I prove the monotonicity Lemma which is based on the standard revealed preferences argument.

**Lemma 6.** *For any $q_L^* \in \operatorname{supp} \sigma_L^*$, $q_H^* \in \operatorname{supp} \sigma_H^*$ holds*

$$C(q_H^*)a^*(q_H^*) \leq C(q_L^*)a^*(q_L^*)$$

*In a PBE[12] such that $a^*(q_H^*) = 1$ holds*

$$q_H^* \leq q_L^*$$

---

[12] $a^*(q_H^*) = 1$ will follow from (32) for "most" of the equilibria. It will follow from Lemma 9 that $q_H^* \leq q_L^*$ even if $a^*(q_H^*) < 1$. We need to study the structure of the equilibria in more details to prove that $q_H^* < q_L^*$ holds in almost any PBE.

Proof is given in the Appendix.

Next, I show that if an equilibrium has a pooling part, it can consist of only one offer.

**Lemma 7.** *If* $\operatorname{supp} \sigma_H^* \cap \operatorname{supp} \sigma_L^* \neq \varnothing$ *then there is only one common point in the supports of the equilibrium mixed strategies for the two types:*

$$\operatorname{supp} \sigma_H^* \cap \operatorname{supp} \sigma_L^* = \left\{ q_p^* \right\}$$

Proof is given in the Appendix.

The two Lemmas 6 and 7 show that any PBE has a very particular structure: there can be a pooling part - an offer $q_p^*$ made by both types of the Principal, and a separating part - the offers made only by $\alpha_H$-type lying to the left of $q_p^*$, and the offers made only by $\alpha_L$-type to the right of $q_p^*$. Put formally,

$$q_H^* < q_p^* < q_L^* \tag{13}$$

for any $q_H^* \in \operatorname{supp} \sigma_H^* \setminus \operatorname{supp} \sigma_L^*$, $q_L^* \in \operatorname{supp} \sigma_L^* \setminus \operatorname{supp} \sigma_H^*$ if the corresponding elements of an equilibrium exist.

I will distinguish between

- pooling equilibria - equilibria with $\operatorname{supp} \sigma_H^* = \operatorname{supp} \sigma_L^*$,

- semi-separating equilibria[13] - equilibria with $\operatorname{supp} \sigma_H^* \neq \operatorname{supp} \sigma_L^*$ and $\operatorname{supp} \sigma_H^* \cap \operatorname{supp} \sigma_L^* \neq \varnothing$

- separating equilibria, for which $\operatorname{supp} \sigma_H^* \cap \operatorname{supp} \sigma_L^* = \varnothing$.

Following lemma characterizes the set of pooling equilibria.

**Lemma 8.** *The pooling equilibria always exists. In any pooling equilibrium there is only one offer* $q_p^*$ *made by both Principal types:* $\operatorname{supp} \sigma_H^* = \operatorname{supp} \sigma_L^* = \{q_p^*\}$.
*The offer* $q_p^*$ *is a pooling equilibrium offer iff* $q_p^* \leq q_E^0$
*Acceptance rule necessary satisfies*[14]

$$a^*(q_p^*) = \begin{cases} 1 & \text{if } q_p^* < q_E^0 \\ any \quad a \in [0,1] & \text{if } q_p^* = q_E^0 \end{cases}$$

*and belief necessary satisfies* $\mu^*(q_p^*) = \Pi$.

Proof is given in the Appendix.

Figure 3 below illustrates the set of pooling equilibria.

I turn now to the separating and semi-separating equilibria.

For the further analysis the relative position of $q_H^P, q_L^P, q_L^0, q_E^0$ is important. It is partially described by (8) and (9). More precise characterization can be obtained by inspecting Figure 2. The formal statement is as follows.

**Proposition 1.** *There exist thresholds* $\beta_3 < \beta_2 < \beta_1$ *determined by*

$$\begin{aligned} \beta_1 & \quad \text{is solution to} \quad q^0(\alpha_L, \beta) = q_L^P \\ \beta_2 & \quad \text{is solution to} \quad q^0(\alpha_L, \beta) = q_H^P \\ \beta_3 & \quad \text{is solution to} \quad q^0(E\alpha, \beta) = q_H^P \end{aligned}$$

---

[13]Semi-separating equilibria are often called "hybrid" in the literature.

[14]Additional requirements should be imposed on the acceptance rule for the offers $q \neq q^*$ to ensure that the two types of Principal don't have profitable deviations.
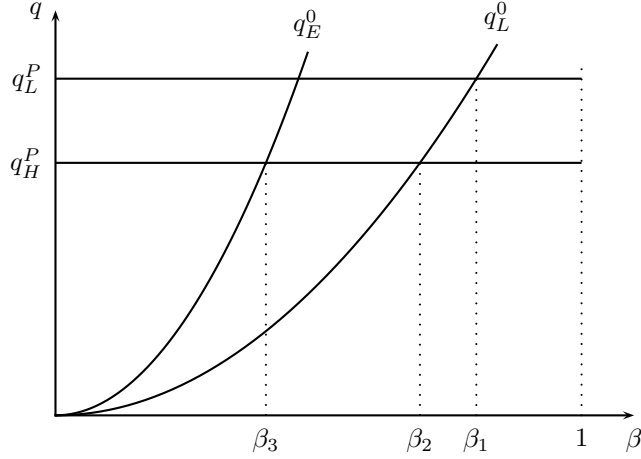
Figure 2: The relative position of $q_H^P, q_L^P, q_L^0, q_E^0$

*such that*

$$q_L^0 < q_H^P \quad \text{for } \beta < \beta_2$$
$$q_H^P < q_L^0 \quad \text{for } \beta > \beta_2$$

The thresholds $\beta_j$ are important to describe the equilibria structure and properties.

For the analysis of the separating and semi-separating equilibria the crucial is the relative position of $q_H^P$ and $q_L^0$, i.e. whether $\beta > \beta_2$ or $\beta < \beta_2$.

The following lemma characterizes the structure of the semi-separating and separating equilibria of the signaling game.

**Lemma 9.** *The separating equilibria always exist.*
*For $\beta \geq \beta_2 \iff q_L^0 \leq q_H^P$ the supports of the equilibrium offers consist of only one offer*

$$\text{supp } \sigma_j^* = \{q_j^*\}, \ j = L, H$$

*and*

$$q_H^* < q_L^* = q_L^0, \quad a^*(q_H^*) = 1, \ a^*(q_L^*) < 1$$

*For $\beta < \beta_2 \iff q_H^P < q_L^0$ the supports of the equilibrium offers can consist of one or two offers.*

*The semi-separating equilibria always exist.*
*For $\beta \leq \beta_2 \iff q_L^0 \leq q_H^P$ they can have one of the two structures:*

$$q_H^* < q_p^*, \quad a^*(q_H^*) = 1, \ a^*(q_p^*) \leq 1$$
$$q_p^* < q_L^* = q_L^0, \quad a^*(q_p^*) = 1, \ a^*(q_L^*) < 1$$

*For $\beta > \beta_2 \iff q_H^P < q_L^0$ they can also have the third structure*

$$q_H^* < q_p^* < q_L^*, \quad a^*(q_H^*) = a^*(q_p^*) = 1, \ a^*(q_L^*) \leq 1$$
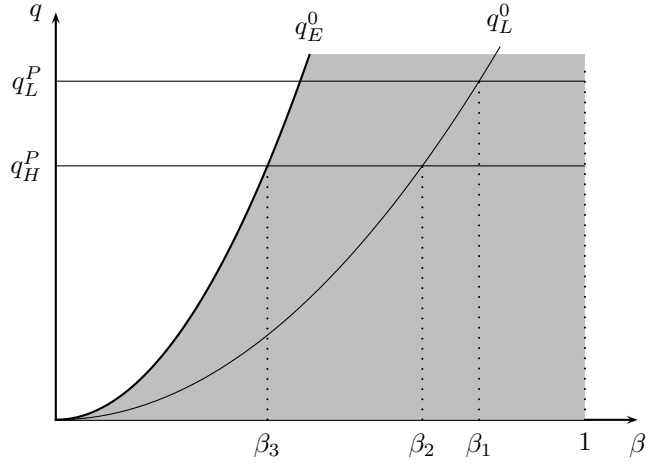
11

Figure 3: Pooling equilibria

Proof is given in the Appendix.

Now we have obtained the structure of the set of PBE of the signaling game in great details. Figure 3 represents the set of all the pooling equilibria. For each $\beta$ any point from the shadow area represents a pooling equilibrium offer.
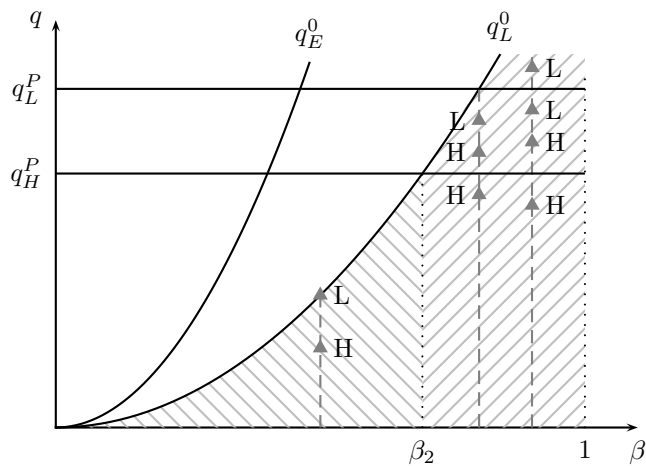


Figure 4: Separating equilibria

Figure 4 represents the set of the separating equilibria and shows some typical cases, as described by Lemma 9.

The figure shows three equilibrium offers profiles corresponding to the three different values of $\beta$. The triangles represents the equilibrium offers and the labels show which type makes which offer.

In the left-most equilibrium, $\alpha_L$-type offers $q_L^* = q_L^0$ which is accepted with probability $a^*\left(q_L^0\right) < 1$. The $\alpha_H$-type makes an offer $q_H^* < q_L^0$ which is ac-

12

cepted with probability 1. If $a^* \left( q_L^0 \right)$ is not too low and $q_H^*$ isn't too small, this configuration is incentive compatible and constitutes an equilibrium.

For the equilibrium in the middle of the figure, $\alpha_H$-type is indifferent between the two offers, one of which is higher and another is lower than $q_H^P$; $\alpha_L$-type makes only one offer. All the equilibrium offers are accepted with probability 1. All the other offers are rejected, according to the acceptance rule[15].

Finally, in the right-most equilibrium both types use mixed offers with 2-points support.
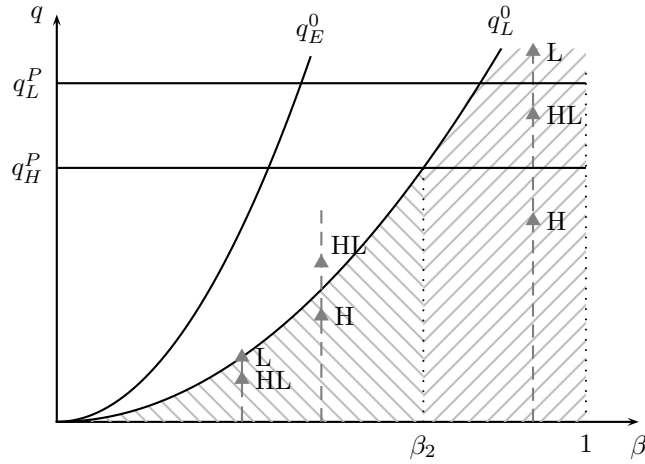


Figure 5: Semiseparating equilibria

Figure 5 represents the set of the semi-separating equilibria.

The multiplicity of equilibria is the consequence of no restrictions on the acceptance rule for the out-of-equilibrium offers. For example, consider the right-most equilibrium of figure 5. Notice that the Principals' preferred offers $q_j^P$ are below $q_L^0$ which means that it is in the Agent's interest to accept these offers would they be proposed. However, since they are out of equilibrium, there is no rationality restrictions for the sequential continuation subgames moves started when these offers are made. Consequently, these offers can be rejected (or accepted with low probability) in the equilibrium acceptance rule. In this way the equilibrium is sustained to make deviation to these offers unprofitable for the Principal. However, such acceptance rule is "unreasonable" and the outcome of the game predicted by this equilibrium can hardly appear.

Equilibrium refinement is a standard procedure for the signaling games which allows to eliminate "unreasonable" equilibria.

---

[15]Of course, if the out-of-equilibrium offers are accepted with low enough probabilities, such profile can also constitute an equilibrium.

### 3.1.2  Equilibrium Refinement

I use the intuitive criterion[16] proposed by Cho and Kreps [1987]. I show that only one PBE satisfies the criterion for $\beta > \beta_2$. The intuitive criterion also eliminates some equilibria for the case $\beta < \beta_2$, though there is still a continuum of equilibria satisfying it. I will argue then that one of them is "the most reasonable" by applying stronger refinement.

Some notation is needed to implement the intuitive criterion, following Cho and Kreps [1987]. Let $\nu(\mathcal{A}'|q)$ be the ex-post probability which the Agent assigns to the (sub)set $\mathcal{A}'$ of the Principal's type after observing an offer $q$. For the two-type case $\mathcal{A}'$ can be $\{\alpha_H\}$, $\{\alpha_L\}$ or $\mathcal{A} = \{\alpha_H, \alpha_L\}$.

Formally,

$$\nu(\mathcal{A}'|q) = \mu(q)I_{\{\alpha_H \in \mathcal{A}'\}} + (1 - \mu(q))I_{\{\alpha_L \in \mathcal{A}'\}}$$

Let

$$BR(\mathcal{A}', q) = \bigcup_{\mu:\ \nu(\mathcal{A}'|q)=1} BR_\mu(q)$$

be the set of all reasonable acceptance rules for beliefs concentrated on the (sub)set $\mathcal{A}'$ applied to the offer $q$ would it be proposed. The acceptance rule is "reasonable" if it is a Best Response corresponding to some beliefs concentrated on the (sub)set $\mathcal{A}'$.

Fix an equilibrium profile $((\sigma_H^*(\cdot), \sigma_L^*(\cdot)), a^*(\cdot); \mu^*(\cdot))$. Denote the equilibrium payoff of $\alpha_j$-type Principal by $V_j^* = V(q_j^*, a^*(q_j^*), \alpha_j)$ for some[17] $q_j^* \in \text{supp } \sigma_j^*$.

Let $J(q)$ be the set of types which for sure don't want to deviate to $q$ for any reasonable acceptance rule:

$$J(q) = \left\{ \alpha_j : \ V_j^* > \max_{a \in BR(\mathcal{A}, q)} V(q, a, \alpha_j) \right\}$$

The set $\mathcal{A} \setminus J(q)$ then is the set of Principal's types which for sure want to deviate from an equilibrium and make an offer $q$ provided that some reasonable acceptance rule will be applied.

The equilibrium satisfies the intuitive criterion if

$$V_j^* \geq \min_{a \in BR(\mathcal{A} \setminus J(q), q)} V(q, a, \alpha_j) \quad \text{for all } j, q \tag{14}$$

The following Proposition states the main result of the formal analysis of the signaling game.

**Proposition 2.** *For the case $\beta > \beta_2$ $\left( \Leftrightarrow q_H^P < q_L^0 \right)$ the only equilibrium outcome satisfying the intuitive criterion is the outcome with*

$$q_H^* = q_H^P, \qquad q_L^* = \min\left\{ q_L^0, q_L^P \right\}$$

*so that the separating equilibrium emerges.*

*For the case of $\beta \leq \beta_2$ $\left( \Leftrightarrow q_L^0 \leq q_H^P \right)$ only the pooling equilibria with*

$$q_L^0 \leq q_p^* \leq q_E^0$$

*satisfy the intuitive criterion.*

---

[16]A survey of the refinements procedures and approaches can be found in Fudenberg and Tirole [1991].

[17]Clearly, for any $q_j^* \in \text{supp } \sigma_j^*$, $V(q_j^*, a^*(q_j^*), \alpha_j)$ takes the same value, so the definition of $V_j^*$ is correct.

Proof is given in the Appendix.

After applying the intuitive criterion we have obtained a unique prediction for the signaling game outcome for $\beta > \beta_2 \left( \Leftrightarrow q_H^P < q_L^0 \right)$. However, for the case of $\beta \leq \beta_2 \left( \Leftrightarrow q_L^0 \leq q_H^P \right)$ there are still many pooling equilibria, though fewer than in the set of PBE. It can be shown that they can't be eliminated by applying Criterion D1 or NWBR criteria.

The reason for multiplicity of the pooling equilibria satisfying the intuitive criterion is that the equilibrium payoff for the Principal is compared in (14) with the worst (for the Principal) reasonable acceptance rule based on Agent's belief concentrated on the set $\mathcal{A} = \{\alpha_H, \alpha_L\}$. In the worst case, after observing a deviation to $q > q_{HL}^*$ the Agent believes that this deviation is done by $\alpha_L$-type. Then, any offer $q > q_L^0$ is reasonably rejected and the intuitive criterion (14) is satisfied for PBE offers $q_p^* > q_L^0$ since any upward deviation from $q_p^*$ is reasonably rejected.

However, once such upward deviation is profitable for both types, the intuitive criterion can be strengthened by comparing the equilibrium payoff in (14) with payoff obtained under "more reasonable" acceptance rule which is based on the ex-ante belief instead of the worst belief. Then, in the right-hand side of (14) the acceptance rule applied to deviations $q > q_{HL}^*$ is the Best Response acceptance rule based on belief $\mu(q) = \Pi$, so that any $q < q_E^0$ (not only $q < q_L^0$) is accepted. This rules out all the PBE with $q < q_E^0$.

For an equilibrium offer $q_p^* = q_E^0$ the Agent is indifferent between accepting and rejecting (and between any probability of accepting), but for the acceptance rule satisfying the strengthened intuitive criterion, all the offers just below $q_E^0$ are accepted with probability 1, so if $a^*(q_E^0) < 1$, then there will be profitable deviation for any Principal to $q_E^0 - \varepsilon$. Consequently, only $a^*(q_E^0) = 1$ is possible in a (refined) equilibrium.

The strengthening of the intuitive criterion in this way is equivalent to requiring the acceptance rule to be sequentially rational on the out-of-equilibrium path. It is also equivalent to eliminating the weakly dominated acceptance rules.

As a result, only equilibrium with the offer $q_p^* = q_E^0$ for $\beta \leq \beta_3$ satisfies the strengthened intuitive criterion. For $\beta_3 < \beta < \beta_2$ the strengthened criterion also leads to the unique prediction for the game outcome $q_p^* = q_H^P$. Indeed, this offer is feasible, and $\alpha_H$-type prefers to make it. By deviating to a higher offer, $\alpha_L$-type would be revealed and then the deviating offer would be rejected, so $\alpha_L$-type has to pool on $q_H^P$.

It's easy to check that the unique separating equilibrium for $\beta > \beta_2$ which satisfies the intuitive criterion satisfies the strengthened criterion as well.

## 3.2 The 2-type Signaling Game Outcome

The intuitive criterion gives a unique prediction of the game outcome for the case of separating equilibrium, which emerges for $\beta > \beta_2$. For $\beta \leq \beta_2$ any equilibrium satisfying the intuitive criterion is pooling, and the unique prediction for the outcome is obtained with the strengthened intuitive criterion, which selects the Pareto-dominating outcome. The unique (refined) equilibrium outcome is presented by Figure Figure 6 for the case when all the thresholds $\beta_j$ are less than 1. When some of the thresholds $\beta_j$ ar greater than 1, they should be cut to 1.
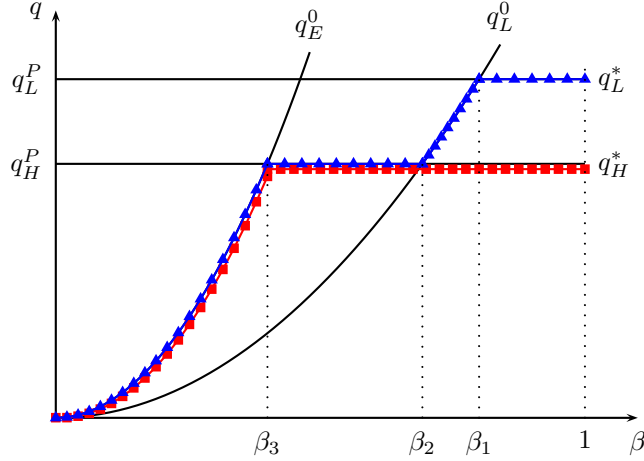
Figure 6: Equilibrium in the signaling game

For the rest of the discussion I call $\alpha_L$-type the "tough" Principal and $\alpha_H$-type - the "generous" Principal. Consider an intuitive explanation for the signaling game outcome.

When the Agent is highly reciprocal, $\beta \in (\beta_1, 1]$, he can be intrinsically motivated even by tough Principal and agrees to perform at her preferred level $q_L^P$, which is quite high, because the Agent's participation constraint isn't very tight. In this case, the tough Principal will reveal her type by requiring performance at this high level.

However, as the Agent's reciprocity declines, his intrinsic motivation decreases, and the tough Principal can't inspire the Agent to perform at the level $q_L^P$ and has to follow the Agent participation constraint by offering contract with $q < q_L^P$, still revealing her type for $\beta \in (\beta_2, \beta_1)$. The Agent gets zero social utility in this case.

For even lower reciprocity intensity, the Agent's intrinsic motivation is not enough to make him exert effort higher than $q_H^P$ if he would learn that the Principal is tough. In this case, the tough Principal follows the offer of the generous Principal, in other words, she has to mimic the generous Principal. The Agent can't distinguish the two types of the Principal and his participation threshold becomes being based on the expected value of the Principal altruism $E\alpha > \alpha_L$ so that $q = q_H^P$ breaks even. The generous Principal doesn't want to separate since she gets her preferred output $q_H^P$. The Agent's expected utility is positive. This is the case for $\beta \in (\beta_3, \beta_2)$.

Finally, when the Agent's reciprocity intensity is very low so that the Agent doesn't have enough intrinsic motivation to exert effort $q_H^P$, both types of Principal follow the Agent's participation threshold for the expected value of the Principal's altruism $E\alpha$, not revealing the type. Neither type of the Principal has an incentive to deviate and reveal his type. This is the case for $\beta \in (0, \beta_3)$.

To conclude the description of the outcome, I stress that in the 2-type signaling game the tough Principal doesn't want to mimic the generous one unless she has to, because mimicking will result in lower output which is not desirable

16

for the tough Principal. However, the Agent is (intrinsically) motivated to accept an offer if it isn't too high. As a result, if the Agent's intrinsic motivation isn't high enough - due to low intensity of reciprocity or due to revealed low altruism of the tough Principal, the tough Principal has to follow the Agent's participation threshold or the generous Principal's offer.

## 3.3 Signaling with Continuum of Types

The analysis for 2 types can be generalized to the case of a continuum of types. I don't present the complete analysis, as for the 2-type case and rather focus on the "most reasonable" equilibrium[18], which leads to the unique prediction of the signaling game outcome.

Let the Principal's altruism parameter $\alpha$ be distributed on the interval $[\alpha_1, \alpha_2] \subset [0, 1]$ with continuous CDF $F(\alpha)$.

The interval's bounds $\alpha_1$ and $\alpha_2$ are the exact bounds of the distribution: $\alpha_1 = \inf\{\alpha | F(\alpha) > 0\}$, $\alpha_2 = \sup\{\alpha | F(\alpha) < 1\}$.

Let $\alpha^\times$ be the solution to

$$q^0\left(\alpha^\times, \beta\right) = q^P\left(\alpha^\times\right)$$

which always exists[19] and is unique since the left-hand side is increasing and the right-hand side is decreasing function of $\alpha$ - see Figure 7.

The next property follows directly from the definition of $\alpha^\times$.

**Claim 1.** *The Principal's preferred output $q^P$ is feasible, i.e. satisfies the Agent's participation constraint iff $\alpha \geq \alpha^\times$.*

So, the population of the Principals can be separated into two sub-populations. One consists of comparatively generous ones with $\alpha \geq \alpha^\times$, which can inspire the Agent's intrinsic motivation high enough to implement their preferred output $q^P$, would their altruism be revealed. Another subpopulation consists of comparatively tough Principals with $\alpha < \alpha^\times$, which can not inspire high enough Agent's intrinsic motivation.

Denote by

$$E_{\widetilde{\alpha}}[\alpha] = E[\alpha | \alpha < \widetilde{\alpha}]$$

the expected value of the truncated distribution of Principal's altruism parameter, bounded at the top by $\widetilde{\alpha}$.

The signaling game outcome in the "most reasonable" equilibrium is characterized by the following Proposition.

**Proposition 3.** *1. If $\alpha$ is distributed inside the interval $[\alpha^\times, 1]$, then all the Principals implement their preferred output in the "most reasonable" equilibrium, i.e. the equilibrium contract is*

$$q = q^P(\alpha)$$

*2. If $\alpha_1 < \alpha^\times$ and*

$$q^0\left(E\alpha, \beta\right) > q^P\left(\alpha_2\right) \tag{15}$$

---

[18]See the refinement section for the 2-type case for the discussion of the "most reasonable" equilibrium.

[19]The fact that it can be that $\alpha^\times > 1$ is not a problem.

17

*then there exists $\tilde{\alpha} \in [\alpha^\times, \alpha_2]$ determined as solution to*

$$\tilde{q}^0 \equiv q^0 \left( E_{\tilde{\alpha}}[\alpha], \beta \right) = q^P \left( E_{\tilde{\alpha}}[\alpha] \right) \tag{16}$$

*such that the "most reasonable" equilibrium contract is given by*

$$q = \begin{cases} q^P(\alpha, \beta) & \text{for } \alpha > \tilde{\alpha} \\ \tilde{q}^0 & \text{for } \alpha \le \tilde{\alpha} \end{cases}$$

*where $\tilde{q}^0$ is determined by (16).*

3. *If the inequality (15) doesn't hold[20], then the "most reasonable" equilibrium contract is full pooling with*

$$q = q^0 \left( E\alpha, \beta \right)$$

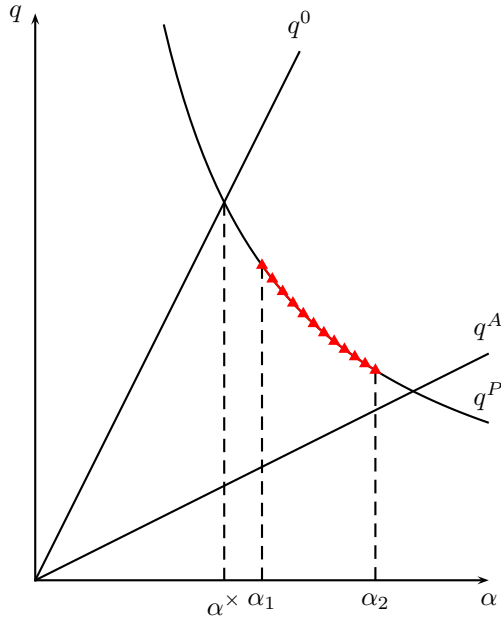Proof is given in the Appendix.



Figure 7: Equilibrium in signaling game. Case 1

Figure 7 illustrates the equilibrium contract for the case when ex-ante all the Principals in the population are highly-altruistic, i.e. $\alpha \ge \alpha^\times$ and can inspire the Agent to perform at their preferred level (point 1 of the Proposition).

Figure 8 illustrates the case when there are both generous and tough Principals, i.e. with altruism greater and smaller than $\alpha^\times$, but there are enough Principals with high altruism (point 2 of the proposition). In this case only a part of the Principals' subpopulation with $\alpha > \alpha^\times$, for instance those with $\alpha \ge \hat{\alpha}$ implement their preferred performance level $q^P$.
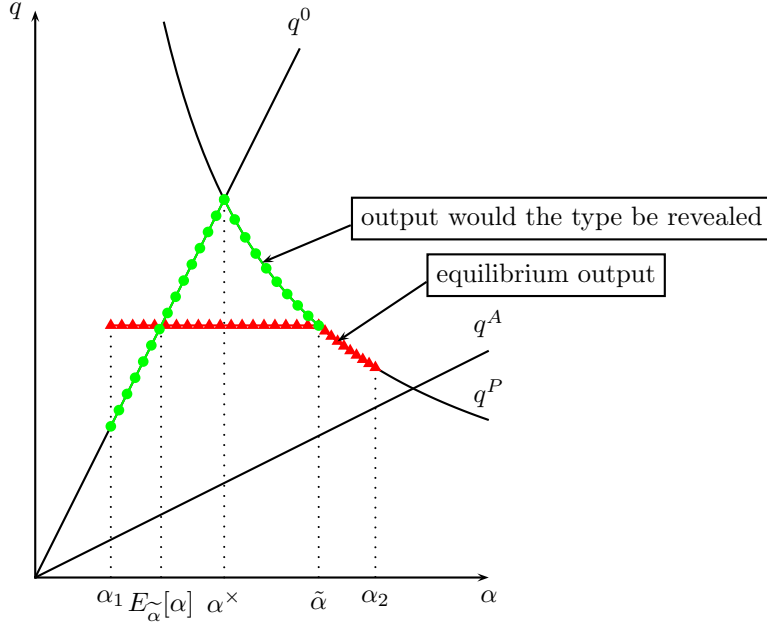
---

[20]This implies $\alpha_1 < \alpha^\times$.

18

Figure 8: Equilibrium in signaling game. Case 2

The first case from Proposition 3 emerges for high $\beta$. As $\beta$ decreases, the equilibrium structure switches to the one described at point 2 of the Proposition and then to the one of point 3.

Similarly to the 2-type case, when the Agent is highly reciprocal, the separating equilibrium in which all the Principals implement their preferred output emerges. For lower levels of reciprocity, the equilibrium structure shifts to pooling and a larger share of Principals can't implement their preferred performance level.

## 3.4  Application to the Organization Design

For the 2-type framework, the generous Principal implements her preferred output for $\beta \in [\beta_3, 1]$. In this case she cares neither about the altruism level of the tough Principals $\alpha_L$, nor about the structure of the Principals' population characterized by $\Pi$.

However, if $\beta \in [0, \beta_3)$, the generous Principal can't implement her preferred output in the emerging pooling equilibrium. In this case she is affected by the adverse effect which emerges due to the very existence of the tough Principals in the population.

Such adverse effect is present for the continuum-type framework as well. In particular, for the case illustrated by Figure 8, the Principals with $\alpha \in [\alpha^\times, \tilde{\alpha})$ are affected. For the case, described in point 3 of Proposition 3, the whole generous subpopulation with $\alpha \in [\alpha^\times, \alpha_2]$ is affected.

Even a part of the "tough" Principals with $\alpha < \alpha^\times$ is affected as they get output lower than the Agent's participation threshold would the Principals' type

be revealed. For the case depicted by Figure 8, these have $\alpha \in (E_{\widehat{\alpha}}[\alpha], \alpha^\times]$; for the case of point 3 of Proposition 3 these have $\alpha \in (E\alpha, \alpha^\times]$.

In both frameworks - the 2-type and the continuum-type,

- the most altruistic Principals (among the generous ones) are less likely to be affected,

- if Agent's reciprocity is high, then all the concerned Principals are less likely to be affected,

- if tough Principals are not too tough or if they represent a smaller share in the population (so that the expected altruism level in the whole population is high enough), then all the concerned Principals are less likely to be affected.

This means that the most generous Principals care to a lesser extent about the environment in which they operate. For the Principals which are generous but not the most generous ones, the environment can play an important role for their ability to elicit a high performance level and they prefer to be in the environment of the generous Principals and highly reciprocal Agents.

The counterpart of the adverse effect is an advantageous effect for the tough Principals which emerges due to the very existence of the generous Principals in the population. In the 2-type framework, for the case of $\beta \in [0, \beta_2)$ the tough Principal gets an output which exceeds the Agent's participation threshold would the type of the tough Principals be revealed. In the continuum-type framework, some Principals benefit from the same advantage: for the case depicted by Figure 8, these are ones with $\alpha \in [\alpha_1, E_{\widetilde{\alpha}}[\alpha])$, for the case depicted by Figure ??, these are ones with $\alpha \in [\alpha_1, E\alpha]$. The advantageous effect appears when the tough Principal can mimic to be the generous ones (or when these more generous can't separate from the tougher ones).

As a result, the tough Principals prefer having highly reciprocal Agents and if this is not the case - then to be in the environment of the generous Principals to benefit from the advantageous effect in the pooling equilibrium.

Now let us take a point of view of an organization designer. Assume that he has a choice - whether to create an organizations of type A in which the Principals' population is mixed, or create an organization of type B in which the generous and tough Principals are separated. The organization designer informs in this way the Agents about the Principal's altruism (the information can still be imprecise).

For instance, in the 2-type framework, the organization designer can separate the two types, in the continuum-type framework, the designer can separate the Principals with $\alpha > \alpha^\times$ from those below the threshold $\alpha^\times$.

An interesting question is how the choice of organization type influences the overall performance in the organization?

Clearly, in organizations of type A the generous Principals will implement lower output, whereas in the organizations of type B the generous Principals will get higher output. This shows that if separation emerges as a result of signaling (endogenously), the tough Principals implement higher output; if separation is exogenous then the relation between the outputs implemented by the tough and generous Principals is reversed.

Figure 9 illustrates exogenous Principals' separation compared to the endogenous Principals' population for the continuum-type framework.
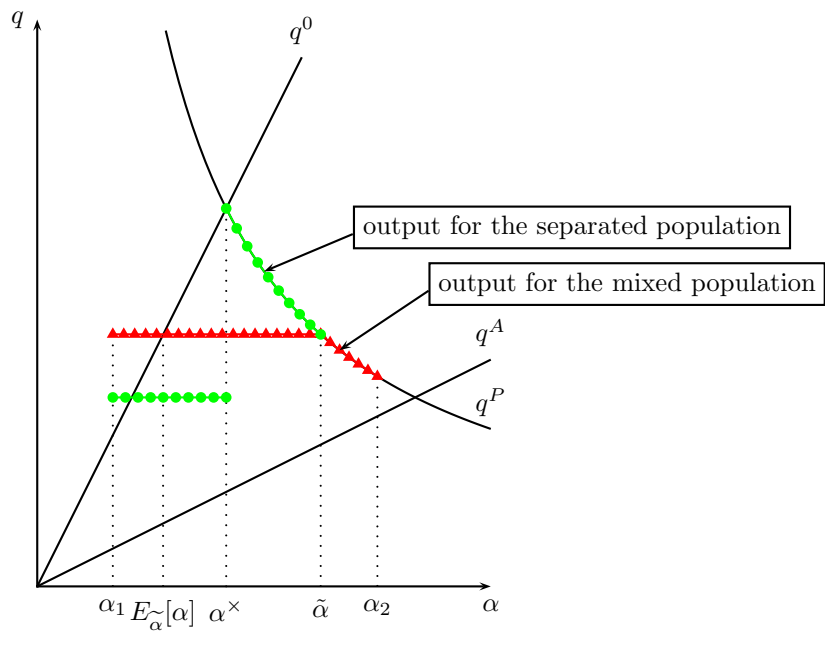
Figure 9: Application to the organization design. The effect of exogenous separation of the Principals' population

# 4    Application to the Experiments

In this section I introduce extrinsic incentives in the Core Model. The modified versions of the Core model are applied to the analysis of human behavior in the settings of two well known lab experiments - the Trust Game of Fehr and Rockenbach [2003] and the Control Game[21] by Falk and Kosfeld [2006], which are described in detail below.

Under the selfishness assumption, which is standard in (neoclassical) contract theory, the results of the experiments are puzzling whereas the reciprocal altruism framework allows to account for the behavior observed in the lab, for instance the motivation crowding out. This also shows the potential of the modeling approach based on reciprocal altruism and provides an additional justification for the relevance of contracting models based on reciprocal altruism.

## 4.1    Trust Game

Consider the Trust Game (or Investment Game) in its Fehr and Rockenbach [2003] version. In their experiment, both the Principal and the Agent are endowed with $S = 10$ units of money. First, the Principal decides on $x$ - how much money to send to the Agent and also announces $\widehat{q}$ - the desired back-transfer. The desired back-transfer isn't binding for the Agent. The experimenter triples the sum of money sent by the Principal[22], so that the Agent receives $3x$. The Agent then decides on the back-transfer $q$. This setting represents the *Baseline treatment*. Notice that $\widehat{q}$ is a "cheap talk" in this case.

In the *Incentive treatment* the Principal on top of $x$ and $\widehat{q}$ announces a fine $f$ which is imposed on the Agent if the back-transfer is lower than the desired level $\widehat{q}$, so that $\widehat{q}$ is no more a "cheap talk". The fine isn't paid to the Principal, it simply reduces the Agent's payoff, in other words the fine is a punishment for the Agent. The fine amount is exogenous (set by the experimenter), so that the only decision of the Principal is to choose whether to impose the fine or not. The contract is now $(x, \widehat{q}, f)$ where $f$ can take only two values 0 or $\overline{f}$.

The paper brings evidence of crowding-out effect, i.e. the situation in which the use of the extrinsic motivator (fine for bad performance) decreases the intrinsic motivation and, as a result, leads to a lower performance. The study finds that, on average, in the incentive treatment the back-payment is higher when the fine is set to zero (the Principal chooses not to punish) than for the case of punishing ($f = \overline{f}$). This means that imposing (extrinsic) incentive leads to a lower performance, which is puzzling from the point of view of "standard" contract theory.

The utilities of the Principal and the Agent in the experimental setting are

---

[21]I follow Ellingsen and Johannesson [2008] in calling the game of Falk and Kosfeld [2006] the "Control Game".

[22]This explains why the game can also be called the "Investment Game". The transfer $x$ can be thought of as an investment, $3x$ - as a return to the investment.

given, following the reciprocal altruism framework[23] by

$$V = 10 - x + q + \alpha(10 + 3x - C(q) - fI_{q<\widehat{q}})$$
$$U = 10 + 3x - C(q) - fI_{q<\widehat{q}} + \widehat{\alpha}\beta(10 - x + q)$$

Let the decision on $x$ has already been made and focus on the continuation subgame[24] in which the Principal decides on $\widehat{q}$ and $f$ and then the Agent decides on $q$. We can consider $x$ as a constant at this point. Dropping $x$ and other constants in the payoff functions leads to the following simplified expressions for the payoffs:

$$V = q - \alpha(C(q) + fI_{q<\widehat{q}}) \tag{17}$$
$$U = \widehat{\alpha}\beta q - C(q) - fI_{q<\widehat{q}} \tag{18}$$

In the absence of a fine, the Principal's preferred back-transfer $q^P(\alpha)$ is determined, as in (3) for the Core Model by

$$C'(q^P) = \frac{1}{\alpha}$$

In the absence of the fine, the back-transfer preferred by the Agent $q^A(\widehat{\alpha}, \beta)$ is determined by

$$C'(q^A) = \widehat{\alpha}\beta$$

as in (4) for the Core model.

Consider the following setting.

Let the Principals and the Agents be heterogenous - some of them are pro-social, others are selfish. I denote the type of the Principal by $\theta^P$, and the type of the Agent by $\theta^A$. For both - the Principals and the Agents, $\theta^j \in \{$Social, Selfish$\}$. The type is the private information.

The pro-social actors are characterized by altruism $\alpha_H$ and reciprocity intensity $\beta_H$, the selfish ones - by the pair $(\alpha_L, \beta_L)$. I assume that

$$\alpha_H > \alpha_L, \quad \beta_H > \beta_L, \quad 0 \le \alpha_j, \beta_j \le 1$$

To simplify the analysis, I assume[25] that $\beta_L = 0$.

Assume that the actors are drawn from the same population, in which the share of the pro-social actors $Prob(\theta = $Social$)$ is unknown. By observing own type, the actors update the belief on the partner's type. I assume that the probability, assigned by the pro-social Principal to be matched with the pro-social Agent $\pi_H$ is higher than the probability $\pi_L$, assigned by the selfish Principal to the same event:

$$\pi_H = Prob(\theta^A = \text{Social}|\theta^P = \text{Social}) = Prob(\beta = \beta_H|\alpha = \alpha_H) \tag{19}$$
$$\pi_L = Prob(\theta^A = \text{Social}|\theta^P = \text{Selfish}) = Prob(\beta = \beta_H|\alpha = \alpha_L) \tag{20}$$

---

[23]In the experiment monetary cost of paying back is linear: $C_m(q) = q$. I assume that there is also a psychological cost of paying back $C_\psi(q)$ which assumed to be convex, so that the overall cost $C(q) = C_m(q) + C_\psi(q)$ is convex. This assumption is admittedly ad hoc, but it is needed to capture the predominance of non bang-bang behavior.

[24]Of course, $x$ itself is a signal of the Principal's altruism, but I assume that the Agent updates his belief on the Principal's altruism after observing $x$ which brings the belief at the beginning of the subgame.

[25]The more general setting with the four possible pairs $(\alpha_k, \beta_l)$ and without requiring $\beta_L = 0$ can be considered. This, however, leads to more tedious computations but doesn't bring additional intuition. So, I restrict attention to the simpler setting.

In other words, I assume that people tend to treat the others as similar to themselves. Such a difference in beliefs can result from the rational projection bias[26].

The Principal moves first and doesn't know the type of the Agent with whom she is matched. The Agent, on the contrary, observes the action of the Principal, and can use this to learn about the Principal's type. Because of this, I assume that behavior of the Principal is driven by her (unconditional) altruism[27], whereas the behavior of the Agent is driven by his reciprocity. In other words, only altruism $\alpha$ affects the Principal's behavior, and only reciprocity intensity $\beta$, interacted with the belief on the Principal's type $\widehat{\alpha}$ - for the Agent's.

This setting brings us to the following signaling game with the two-sided asymmetric information.

The Principal can be of type $\theta^P \in \{$Social, Selfish$\}$, or, equivalently, $\alpha \in \{\alpha_H, \alpha_L\}$. The Agent can be of two types, $\theta^A \in \{$Social, Selfish$\}$, or, equivalently, $\beta \in \{\beta_H, \beta_L\}$.

The Principal's strategy is a type-contingent pair $(f_j, \widehat{q}_j) \in \{0, \overline{f}\} \times [0, +\infty)$, $j = L, H$ (where the index $L$ is used for the selfish type, and $H$ indexes for the pro-social type).

The Agent's strategy is a type-contingent back-transfer conditional on the Principal's action $q_i(f, \widehat{q}) \in [0, +\infty)$, $i = L, H$.

The Principal has ex-post belief on the probabilities to be matched with the pro-social Agent which depends on the Principal's type:

The Agent's ex-post belief $\mu$ is determined by the Principal's observed action, $\mu(f, \widehat{q}) = Prob(\alpha = \alpha_H | f, \widehat{q})$. As for the Core Model, there is a one-to-one correspondence between belief $\mu$ and the ex-post expectation of the Principal's type $\widehat{\alpha}$: $\widehat{\alpha} = \mu \alpha_H + (1 - \mu)\alpha_L$, so that $\widehat{\alpha}$ can be considered instead of $\mu$. The payoffs are given by (17) and (18).

---

[26]Rational projections bias can be determined as "the tendency to look at others (other people or future selves) from the point of view of ones current self" - see Tirole [2002]. For the evidence and implication of the projection bias to "future selves" see, e.g., Loewenstein et al. [2003]. Bénabou and Tirole [2006b] discus the implication of the projection bias for collective belief.

Rational projection bias can be endogeneized in the following way. Assume that the share of the pro-social actors in the population is unknown. It can be either $\Pi_g$ or $\Pi_b$, depending on the state of the world - "Good" or "Bad". Assume that $\Pi_g > \Pi_b$. The state of the world is not observed; however the prior probabilities of the two states are commonly known and equal to $p_g$ and $p_b = 1 - p_g$ respectively.

By observing her own type and using the Bayesian rule, the Principal obtains probabilities to be matched with the pro-social Agent:

$$\pi_H = \Pi_g \frac{p_g \Pi_g}{p_g \Pi_g + p_b \Pi_b} + \Pi_b \frac{p_b \Pi_b}{p_g \Pi_g + p_b \Pi_b} = \frac{p_g \Pi_g^2 + p_b \Pi_b^2}{p_g \Pi_g + p_b \Pi_b}$$

$$\pi_L = \Pi_g \frac{p_g(1 - \Pi_g)}{p_g(1 - \Pi_g) + p_b(1 - \Pi_b)} + \Pi_b \frac{p_b(1 - \Pi_b)}{p_g(1 - \Pi_g) + p_b(1 - \Pi_b)} =$$

$$= \frac{p_g \Pi_g(1 - \Pi_g) + p_b \Pi_b(1 - \Pi_b)}{p_g(1 - \Pi_g) + p_b(1 - \Pi_b)}$$

The objective probability for the Agent to be pro-social is $\pi = Prob(\theta^A = $ Social$) = p_g \Pi_g + p_b \Pi_b$. Clearly, $\pi_L < \pi < \pi_H$.

[27]Alternatively, one can assume that given the prior belief on the Agent's altruism, the Principal's altruism is equal to the sum of her pure (unconditional) altruism $\alpha_p$ and reciprocal altruism $\alpha_r = \beta E[\alpha^A]$. This results in the Principal's altruism towards the Agent at the level $\alpha_H = \alpha_{pH} + \beta_H E[\alpha^A]$ or $\alpha_L = \alpha_{pL} + \beta_L E[\alpha^A]$, depending on the type of the Principal.

I consider the Perfect Bayesian equilibrium, in which Agent's belief off the equilibrium path are "reasonable", as in the intuitive criterion of Cho and Kreps.

There is an important difference with the Core model. Now the choice of the fine $f$ plays crucial role in signaling.

The setting described above corresponds to the Incentive Treatment. For the Baseline Treatment the fine $f$ is exogenously set to zero.

I proceed now backward in the analysis of the game.

Consider the Agent's Best Response back-transfer for the different treatments of the experiment. The Agent's participation threshold isn't relevant since paying back zero is feasible.

**Claim 2.** *In the Trust Game, if the Agent holds beliefs $\widehat{\alpha}$ on the Principal's altruism, the Best Response back-transfer $q$ is:*

1. *in the baseline treatment (fine isn't possible) and in the incentive treatment when the Principal chooses not to punish ($f = 0$): $q = q^A(\widehat{\alpha}, \beta)$.*

2. *in the incentive treatment when the Principal chooses to impose a fine ($f = \overline{f}$):*
$$q = \begin{cases} q^A(\widehat{\alpha}, \beta) & if \ \ \widehat{q} < q^A(\widehat{\alpha}, \beta) \\ \widehat{q} & if \ \ q^A(\widehat{\alpha}, \beta) < \widehat{q} < \widetilde{q}^A(\widehat{\alpha}, \beta) \\ q^A(\widehat{\alpha}) & if \ \ \widehat{q} > \widetilde{q}^A(\widehat{\alpha}, \beta) \end{cases}$$

*where $\widetilde{q}^A(\widehat{\alpha}, \beta)$ is determined by (see figure 10)*

$$\widehat{\alpha}\beta q^A - C(q^A) - f = \widehat{\alpha}\beta\widetilde{q}^A - C(\widetilde{q}^A), \quad \widetilde{q}^A > q^A$$

*and $\widetilde{q}^A(\widehat{\alpha}, \beta)$ is an increasing function of $\widehat{\alpha}$.*
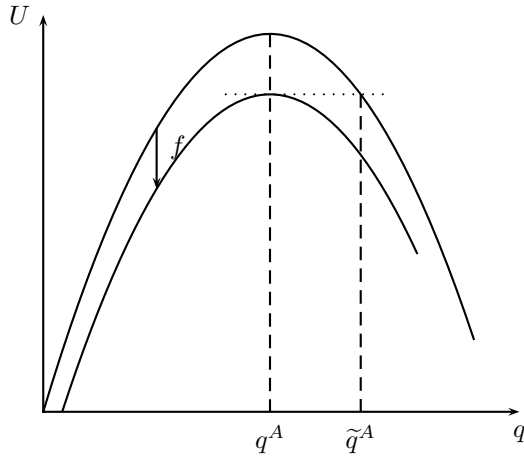


Figure 10: Agent's payoff

Proof is given in the Appendix.

It follows from the Claim that, contrary to the predictions of the standard contract theory, when extrinsic incentives are used for the intrinsically motivated
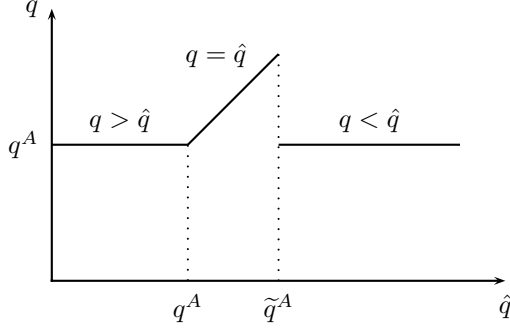
Figure 11: Back-transfer as a function of threshold

Agent, the actual back transfer can be higher, equal or lower than the desirable back-transfer, as illustrated by Figure 11.

If $\widehat{q} < q^A$, the required performance is low compared to the intrinsic motivation, so the Agent is willing to perform better than he is asked for. For $\widehat{q} = q^A$ the intrinsic motivation is just enough to motivate the Agent for the required level of performance. Finally, $\widehat{q} > q^A$ corresponds to the case when intrinsic motivation isn't enough to inspire the Agent for high enough performance.

Not surprisingly, for a given belief $\widehat{\alpha}$ holds $\tilde{q}^A(\widehat{\alpha}) > q^A(\widehat{\alpha})$, which means that under symmetric (or revealed) information, if extrinsic incentive is added to the intrinsic motivation, the performance level is higher. So, the fine serves as a "positive reinforcer" in the sense of Bénabou and Tirole [2003].

When the threat of the fine is imposed, the Agent has to give up some utility and faces a trade-off: whether to diverge from his preferred back-payment $q^A$ to the desired back-payment $\widehat{q}$ or to pay a fine but not diverging from his preferred back-payment $q^A$. The value $\tilde{q}^A$ separates the Agent's preferred option: if the desired back-payment is below this value, the Agent prefers to diverge from $q^A$ to $\widehat{q}$, if the desired back-transfer is higher than $\tilde{q}^A$, the Agent prefers to disobey and pay the fine. So, $q > \tilde{q}^A$ can't be implemented with a threat of fine.

Denote by

$$q_{ij} = q^A(\alpha_i, \beta_j)$$

the maximal back-transfer which can be implemented by using intrinsic motivation only, given that the Agent with $\beta = \beta_j$ holds belief $\widehat{\alpha} = \alpha_i$. These back-transfers are determined by $C'(q_{ij}) = \alpha_i \beta_j$. Since $\beta_L = 0$,

$$q_{HL} = q_{LL} = 0$$

Denote by

$$\tilde{q}_{ij} = \tilde{q}^A(\alpha_i, \beta_j)$$

the maximal back-transfer which can be implemented by using both intrinsic and extrinsic motivation, i.e. by imposing the (threat of) fine, given that the Agent with $\beta = \beta_j$ holds belief $\widehat{\alpha} = \alpha_i$. It follows from Claim 2 that

$$C(\tilde{q}_{LL}) = f \qquad (21)$$

Now I show that under some restrictions on the parameters, the game has the separating equilibrium with crowding-out.

26

**Proposition 4.** *There is a range of parameters* $(\alpha_L, \alpha_H, \beta_H, f, \pi_L, \pi_H)$ *for which there exists unique separating equilibrium satisfying the Cho-Kreps intuitive criterion. This equilibrium is characterized by crowding-out in the Agent's performance.*

*In particular, the parameters should be such that*

$$q_{LH} \leq \widetilde{q}_{LL}, \quad \widetilde{q}_{LH} \leq q_{HH}, \quad \pi_H \geq \widehat{\pi}_H, \quad \pi_L \leq \widehat{\pi}_L$$

*where*

$$\widehat{\pi}_H = \frac{\widetilde{q}_{LL} - \alpha_H C(\widetilde{q}_{LL})}{q_{HH} - \alpha_H C(q_{HH})} \leq 1 \tag{22}$$

$$\widehat{\pi}_L = \frac{\widetilde{q}_{LL} - \alpha_L C(\widetilde{q}_{LL})}{q_{HH} - \alpha_L C(q_{HH})} \geq 0 \tag{23}$$

*In the equilibrium* $\alpha_H$*-type imposes no (threat of) fine,* $\alpha_L$*-type imposes (threat of) fine:*

$$f_L^* = \overline{f}, \quad f_H^* = 0$$

*The desired back-transfers are:* $\widehat{q}_L^* = \widetilde{q}_{LL}$*, any* $\widehat{q}_H^* \leq q_{HH}$*.*
*In particular, it is possible to have*

$$\widehat{q}_L^* > \widehat{q}_H^*$$

*The actual back-transfers are*

$$q_{LL}^* = \widetilde{q}_{LL}, \quad q_{LH}^* = \widetilde{q}_{LL}, \quad q_p^* = 0, \quad q_{HH}^* = q_{HH}$$

*where* $q_{ij}^*$ *is the equilibrium back-transfer to the* $\alpha_i$*-Principal from* $\beta_j$*-Agent.*
*The average back-transfer to the pro-social Principal is higher than that to the selfish one:*

$$\pi q_{HH}^* \geq q_{LL}^*$$

Proof is given in the Appendix.

The proposition states that in the separating equilibrium, the pro-social Principal chooses not to impose the fine, i.e. not to use the extrinsic incentive, signaling in this way her generosity and inspiring high intrinsic motivation of the pro-social Agents. They perform at a comparatively high level $q_{HH}$. If the pro-social Principal meets the selfish Agent, which can't be intrinsically motivated, the performance will be zero: $q_{LL} = 0$, since there are no extrinsic incentives.

Imposing the fine, i.e. using the extrinsic incentive, crowds out the intrinsic motivation of the pro-social Agent because this signals low altruism of the Principal. However, provision of the extrinsic incentive guarantees performance of all Agents[28] at a comparatively low level $\widetilde{q}_{LL}$. This can lead to the observed crowding out in the Agent's performance: the average for the lottery between $q_{HH}$ and $q_{LL} = 0$ is higher than $\widetilde{q}_{LL}$ for some values of the parameters.

The rational projection bias is crucial to guarantee the no-mimicking (incentive compatibility) condition. The Principals have different beliefs on the Agent's type, for instance, the pro-social Principal is more optimistic. Because of this, the pro-social Principal prefers the lottery between high $q_{HH}$ and zero,

---

[28] The pro-social Agents could perform at the level $q_{LH}$, if $q_{LH} > \widetilde{q_{LL}}$. This, however, ruled out by the Proposition assumption.

whereas the selfish Principal prefers the comparatively low but sure outcome[29] $\widetilde{q}_{LL}$.

Let us now discuss the conditions for the separating crowding out equilibrium to emerge, described in proposition 4. The condition $q_{LH} \leq \widetilde{q}_{LL}$ is equivalent to $C(q_{LH}) \leq C(\widetilde{q}_{LL})$. Since $C(\widetilde{q}_{LL}) = \overline{f}$ - see (21), this means that $\overline{f}$ should be high enough:

$$\overline{f} \geq f_1 \equiv C(q_{LH})$$

The back-transfer $\widetilde{q}_{LH}$ is determined, according to Claim 2 by

$$\alpha_L \beta_H q_{LH} - C(q_{LH}) - \overline{f} = \alpha_L \beta_H \widetilde{q}_{LH} - C(\widetilde{q}_{LH})$$

where $\widetilde{q}_{LH}$ is chosen on the decreasing part of the graph of the function $F(q) = \alpha_L \beta_H q - C(q)$ (see Figure 10). Consequently, $\widetilde{q}_{LH} \leq q_{HH}$ is equivalent to

$$\alpha_L \beta_H q_{LH} - C(q_{LH}) - \overline{f} \geq \alpha_L \beta_H q_{HH} - C(q_{HH})$$

which can be rewritten as

$$\overline{f} \leq f_2 \equiv (\alpha_L \beta_H q_{LH} - C(q_{LH})) - (\alpha_L \beta_H q_{HH} - C(q_{HH}))$$

So, $\overline{f}$ shouldn't be very high.

The conditions $\pi_H \geq \widehat{\pi}_H$ and $\pi_L \leq \widehat{\pi}_L$ with the thresholds given by (22) and (23), show that the projection bias should be large enough to ensure no-mimicking.

On the other hand, if $\pi_L$ and $\pi_H$ are fixed, the threshold conditions $\pi_H \geq \widehat{\pi}_H$ and $\pi_L \leq \widehat{\pi}_L$ can require additional restrictions on $\overline{f}$.

For instance, if $\overline{f} = 0$, then $\widetilde{q}_{LL} = 0$, and $\widehat{\pi}_{LL} = 0$. So, condition (23) imposes some restriction $\overline{f} \geq f_1'$. Similarly, (23) imposes some restriction $\overline{f} \leq f_2'$.

The following corollary from the proposition 4, providing a more constructive description of the parameters under which the crowding-out separating equilibrium emerges, can be established.

**Corollary 2.** *For a generic triple $(\alpha_L, \alpha_H, \beta_H)$ and quadratic cost function $C(q) = \frac{c}{2}q^2$ there always exists a non-empty set of the parameters $(\pi_L, \pi_H, \overline{f})$ such that there exists a separating equilibrium with crowding-out.*

*These parameters satisfy to (22) and (23) and the inequality $f_1 \leq \overline{f} \leq f_2$. If cost isn't quadratic, the additional condition*

$$\alpha_L \alpha_H \beta_H^2 \left( \frac{q_{LH}}{\alpha_L \beta_H} - \frac{q_{HH}}{\alpha_H \beta_H} \right) \leq C(q_{LH})$$

*is required.*

Proof is given in the Appendix.

---

[29] This reasoning can also be applied if the Principals and the Agents are drawn from the two independent distribution, so that there is no projection bias, but the Principals are risk-averse. So, the separating equilibrium with crowding-out can also emerge in such setting.

## 4.2 The Control Game

In the experiment conducted in Falk and Kosfeld [2006] the Principal chooses whether to restrict the set of Agent's effort from below. Output is assumed to be equal to effort.

Put formally, the Principal offers a contract $\underline{q}$ which can take two values - 0 or $\overline{q}$, where the latter is exogenously set by the experimenter. The Agent then chooses effort $q \in [\underline{q}, \infty)$. Effort is costly for the Agent. The Agent has an initial endowment of 120.

The experiment has a number of findings which can not be explained within the selfishness framework. For instance, the Agents, proposed a contract $\underline{q} > 0$, exert, on average, less effort, than those with $\underline{q} = 0$, which means that extrinsic incentive (control) has a negative impact on Agents' performance. This represents the hidden cost of control effect. The observed behavior of the Agents was heterogenous: there was observed positive, negative and neutral reaction to control. Finally, many Principals choose not to control.

The reciprocal altruism framework accounts for these experimental findings.

As for the Trust Game, I alter the Core Model to match the experiment design.

The selfish utilities of the Principal and the Agent are given by[30]

$$v = q$$
$$u = 120 - C(q)$$

The reciprocal altruism framework leads to the (social) utilities

$$V = q + \alpha(120 - C(q))$$
$$U = 120 - C(q) + \widehat{\alpha}\beta q$$

The initial endowment of the Agent allows to disregard the Agent's participation constraint. By dropping the constants, the Principal and Agent utilities can be simplified to

$$V = q - \alpha C(q) \tag{24}$$
$$U = \widehat{\alpha}\beta q - C(q) \tag{25}$$

Denote by $q^A(\widehat{\alpha}, \beta)$ the Agent's preferred effort level:

$$q^A = \arg\max_q [U(q; \alpha)] = \arg\max_q [\widehat{\alpha}\beta q - C(q)]$$
$$C'(q^A) = \widehat{\alpha}\beta$$

as for the core model.

Consider the setting with heterogenous Principals and Agents, adopted in the analysis of the Trust Game in subsection 4.1.

The Principal's strategy is a type-contingent choice of control $\underline{q}_i \in \{0, \overline{q}\}$, $i = L, H$. The Agent's strategy is a type-contingent effort, conditional on the Principal's action $q_j(\underline{q}) \in [\underline{q}, +\infty)$, $j = L, H$.

---

[30]The experiment sets $C(q) = q/2$. As for the Trust Game, I assume that $C(q)$ is convex. See footnote 23 for the justification of the assumption.

The Principal's ex-post belief on the probabilities to be matched with the pro-social Agent depends on the Principal's type and are given by (19) and (20). The Agent's ex-post belief is determined by the Principal's observed action, $\mu(q) = Prob(\alpha = \alpha_H | q)$. As for the Core Model, there is a one-to-one correspondence between belief $\mu$ and the ex-post expectation of the Principal's type $\widehat{\alpha}$: $\widehat{\alpha} = \mu\alpha_H + (1 - \mu)\alpha_L$, so that $\widehat{\alpha}$ can be considered instead of $\mu$. The payoffs are given by (24) and (25).

As in the analysis of the Trust game, I consider the Perfect Bayesian equilibrium, in which Agent's belief off the equilibrium path are "reasonable" in the sense of the intuitive criterion of Cho and Kreps.

I proceed backwards in the analysis of the game.

Consider first the Agent's Best Response choice of effort.

**Claim 3.** *If $q^A(\widehat{\alpha}, \beta) \geq \underline{q}$ then the Agent's Best Response is $q = q^A(\widehat{\alpha}, \beta)$; otherwise it is $q = \underline{q}$.*

The Claim is evident as it simply says that the Agent chooses the global maximizer of his utility whenever it's feasible. Otherwise, he chooses the closest feasible effort which is equal to $\underline{q}$ - the lower bound of the set of feasible efforts.

Denote by $q_{ij}$ the effort, voluntarily exerted by the Agent with $\beta_j$ which beliefs that the Principal's type is $\alpha_i$, i.e. $q_{ij} = q^A(\alpha_i, \beta_j)$ and $C'(q_{ij}) = \alpha_i\beta_j$.

Consider now the Principal's decision. If the Principal with altruism $\alpha_i$ holds belief $\pi_i$ to be matched with the pro-social Agent, which, in turn, holds the true belief on the Principal's type, then under no-control her utility is[31]

$$V = \pi_i(q_{iH} - \alpha_i C(q_{iH})) \tag{26}$$

Under control, if effort $q_{iH}$ is available, i.e. $q_{iH} \geq \overline{q}$, and the Principal's utility is still given by (26). If $q_{iH} < \overline{q}$, then the effort $q = \overline{q}$ is implemented and the Principal's utility is

$$V = (\overline{q} - \alpha_i C(\overline{q})) \tag{27}$$

Denote by $q_i^{C1} < q_i^{C2}$ the two roots of the equation (see Figure 12)

$$\pi_i(q_{iH} - \alpha_i C(q_{iH})) = (\overline{q} - \alpha_i C(\overline{q}))$$

Comparing the Principal's utility for the case when the Agent holds the true belief on the Principal's type, given by (26) and (27), leads to the following characterization of the Principal's Best Response in this case:

**Claim 4.** *If the Agent holds the true belief on the Principal's type, then the Principal's Best Response is:*

$$\underline{q} = \begin{cases} \overline{q} & if \ \overline{q} \leq q_i^{C2} \\ 0 & if \ \overline{q} \geq q_i^{C2} \end{cases}$$

Consider now the case when the Principal's type is her private information.
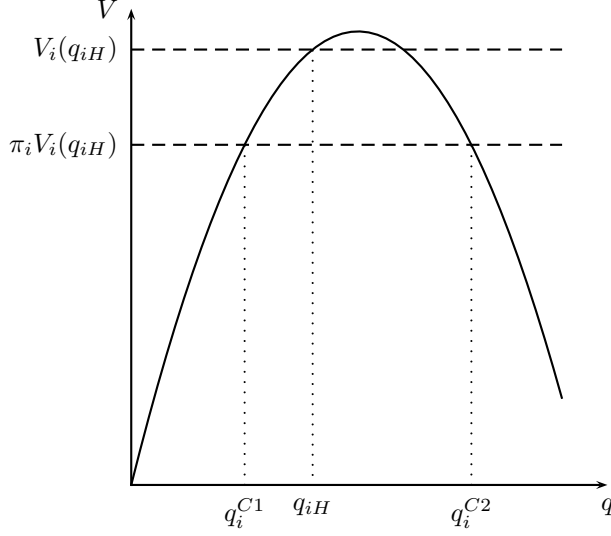
---

[31]$q_{iL} = 0$ since $\beta_L = 0$.

Figure 12: Principal's utility

**Proposition 5.** *For each $\overline{q} < q_{HH}$, $(\alpha_L, \alpha_H, \beta_H)$ there exists a range of parameters $(\pi_L, \pi_H)$ such that there exists the unique separating equilibrium of the Control Game, satisfying the intuitive criterion. There is crowding-out in effort in the equilibrium.*

*The parameters should be such that*

$$\pi_L \leq \widehat{\pi}_L, \qquad \pi_H \geq \widehat{\pi}_H$$

*where for the case $\overline{q} \geq q_{LH}$*

$$\widehat{\pi}_L = \frac{\overline{q} - \alpha_L C(\overline{q})}{q_{HH} - \alpha_L C(q_{HH})} > 0$$

$$\widehat{\pi}_H = \frac{\overline{q} - \alpha_H C(\overline{q})}{q_{HH} - \alpha_H C(q_{HH})} < 1$$

*for the case $\overline{q} < q_{LH}$*

$$\widehat{\pi}_L = \frac{\overline{q} - \alpha_L C(\overline{q})}{[q_{HH} - \alpha_L C(q_{HH})] + [\overline{q} - \alpha_L C(\overline{q})] - [q_{LH} - \alpha_L C(q_{LH})]} > 0$$

$$\widehat{\pi}_H = \frac{\overline{q} - \alpha_H C(\overline{q})}{[q_{HH} - \alpha_H C(q_{HH})] + [\overline{q} - \alpha_H C(\overline{q})] - [q_{LH} - \alpha_H C(q_{LH})]} < 1$$

*In the equilibrium, the pro-social Principal doesn't control, whereas the selfish Principal does:*

$$\underline{q}_H^* = 0, \qquad \underline{q}_L^* = \overline{q}$$

*The Agent's performance*

$$q_{HH}^* = q_{HH}, \qquad q_p^* = 0, \qquad q_{LH}^* = \max\{q_{LH}, \overline{q}\}, \qquad q_{LL}^* = \overline{q}$$

31

where $q_{ij}^*$ is the equilibrium back-transfer to the $\alpha_i$-Principal from the $\beta_j$-Agent.

The average effort to the pro-social Principal is higher than that to the selfish one.

Proof is given in the Appendix.

The mechanism of emerging of the separating equilibrium with crowding-out is similar to that of the Trust Game. By choosing not to control, the pro-social Principal signals her kindness, and this inspires high intrinsic motivation for the pro-social Agent. Because of this, when matched with the pro-social Principal, the pro-social Agent exerts high effort $q_{HH}$. However, the selfish Agent doesn't react to the signal of the Principal's generosity and, once not controlled, exerts zero effort.

The selfish Principal chooses to control and guarantees the (comparatively low) output $\overline{q}$. However, even selfish Principal inspires the pro-social Agent's intrinsic motivation, so effort from the pro-social Agent can be $q_{LH}$, if it's higher than the controlling threshold $\overline{q}$.

The lottery between $q_{HH}$ and 0 effort is differently treated by the two Principals. The pro-social one is more optimistic, and beliefs that the chance to get $q_{HH}$ is higher, compared to the belief of the selfish Principal. Because of this, no-mimicking holds - the pro-social Principal prefers lottery, whereas the selfish Principal prefers the sure outcome.

As for the Trust Game, the separating crowding-out equilibrium emerges when the available extrinsic incentive isn't too weak nor too strong.

I turn now to the more detailed description of other equilibrium structures of the Control Game.

**Proposition 6.** *For given $\alpha_L, \alpha_H, \beta_H, \pi_L, \pi_H$, there exist the threshold values $q_i$, $q_i < q_j$ for $i < j$, such that the equilibrium in the Control Game is:*

1. *No-control pooling for $\overline{q} \in [0, q_1]$, which represents no effect of control;*

2. *Separating equilibrium with crowding-out for $\overline{q} \in [q_2, q_3]$ (hidden cost of control);*

3. *Control pooling for $\overline{q} \in [q_3, q_4]$ (positive effect of control);*

4. *Separating with no crowding-out in effort $\overline{q} \in [q_4, q_5]$ (positive effect of control);*

5. *No-control pooling $\overline{q} \in [q_6, +\infty)$.*

*For $\overline{q} \in [q_1, q_2]$ and $\overline{q} \in [q_5, q_6]$ an equilibrium involves mixed strategies.*

Proof is given in the Appendix.

Figure 13 illustrates the proposition.

The experiment finds that for small $\overline{q}$ the hidden cost of control effect is stronger compared to large $\overline{q}$. The model is in line with this result. In fact, the hidden cost of control is obtained for $\overline{q} \in [q_2, q_3]$; for larger $\overline{q}$, the pooling equilibrium with control emerges which means that the increase in $\overline{q}$ will lead to the increase in average performance (which is equal to $\overline{q}$). For even larger $\overline{q}$ the separating equilibrium in which the controlling Principal gets larger output than the non-controlling one, so there is also the positive effect of control rather than the hidden cost.
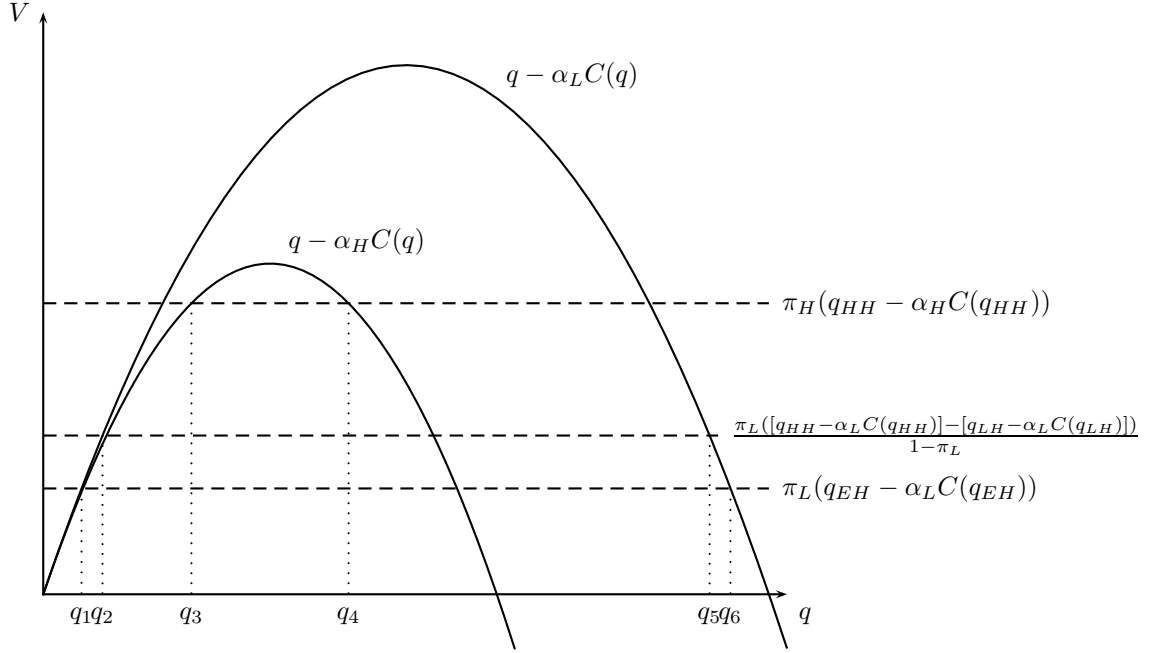
Figure 13: Equilibrium structure in the Control Game

The experiment finds strong heterogeneity between the Agents and that there are Agents which react positively, neutrally or negatively to control. The heterogeneity of the Agents is assumed in the model, and the reaction of the Agents is predicted. Controlling the pro-social Agents can lead to the decrease in his performance. For instance, for $\overline{q} \in [q_2, q_3]$ when the crowding-out equilibrium emerges, the controlled pro-social Agent performs at $q = \overline{q}$, whereas the uncontrolled at $q = q_{HH} > \overline{q}$. Controlling the selfish Agent leads to the increase in his performance for $\overline{q} \in [q_2, q_5]$. Controlling the pro-social Agents also leads to the increase in their performance for the control pooling equilibria. Finally, the neutral effect of control can emerge in the game but out of equilibrium path. According to Claim 3, if Agent's preferred output $q^A$ is higher than control, the Agent will perform at level $q^A$ independently on whether he is controlled or not. In particular, for the case of weak control ($\overline{q}$ is small enough), even if the pro-social Agent is controlled, he can choose to perform at a higher level.

The experiment finds that for larger $\overline{q}$, the larger share of the Principals chooses to control. This is also the prediction of the model: for small $\overline{q}$ none of the Principals controls which results in the no-control pooling equilibrium; for larger values of $\overline{q}$ the hidden cost of control equilibrium emerges in which only the selfish Principals control; for even larger values of $\overline{q}$ the control pooling equilibrium emerges in which both types of Principal choose to control.

## 4.3   A Unifying Framework

The three models considered above - the Core model and its two modifications with a possibility to use an extrinsic motivator have a lot in common and can

be unified in the following way.

First, all of them represent a signaling game.

At stage 1, the Principal sends a message - fixes $\widehat{Q}$ - the set of feasible performance levels for the Agent. Performance is costly for the Agent. At stage 2, the Agent decides whether to comply or disobey. In case of compliance, the Agent chooses an element $q \in \widehat{Q}$ - the implemented performance level. In case of disobeying, the Agent chooses a disobedience option in a set of disobedience options[32] $\overset{\circ}{Q}$: $\overset{\circ}{q} \in \overset{\circ}{Q}$.

For the Core Model $\widehat{Q} = \{q\}$ - the required output level, compliance corresponds to agreement to perform at level $q$ ($a(q) = 1$), disobeying corresponds to $\overset{\circ}{q} = 0$ or, equivalently, to $a(q) = 0$. For the Trust Game the performance set is $\widehat{Q} = [\widehat{q}, +\infty)$, compliance corresponds to the choice of some back-transfer $q \in [\widehat{q}, +\infty)$, disobeying corresponds to the choice of $\overset{\circ}{q} < \widehat{q}$, in the equilibrium with disobedience $\overset{\circ}{q} = q^A$. Finally, for the Control game $\widehat{Q} = (\overline{q}, +\infty)$, compliance means the choice of effort higher than the minimal requirement, disobeying is the choice of effort at the minimal level: $\overset{\circ}{q} = \overline{q}$.

Second, all of the three models are based on the reciprocal altruism framework. This means that the Agent is intrinsically motivated but there is a gap between the Principal's preferred performance level and that of the Agent. The Agent prefers the lower performance.

The weaker the Principal's altruism (i.e. the less generous she is), the larger is the gap, the more the Principal suffers from leaving the Agent too much flexibility, i.e. not restricting at all or imposing only slight restrictions on the performance set $\widehat{Q}$. Consequently, it's in the Principal's interest (especially, the tough Principal) to make a strict offer - i.e. to restrict the Agent's performance set $\widehat{Q}$ to avoid a large deviation from her preferred performance level. On the other hand, offering a slightly restricted or non restricted at all, compared to $\overset{\circ}{Q}$, performance set $\widehat{Q}$ is a generous offer, which is likely to be made by the highly altruistic Principal.

By offering a generous performance set $\widehat{Q}$, the Principal can signal her generosity. By learning the Principal's generosity, the Agent becomes inspired for better performance, i.e. intrinsically motivated. An extrinsic motivator, added to the intrinsic motivation, leads to even better performance for the symmetric information on the Principal's generosity. However, for the case of asymmetric information, the use of an extrinsic motivator signals the non-generosity of the Principal (for the case of separating equilibrium) and crowds out the Agent's intrinsic motivation.

Such a crowding mechanism influences the Agent's performance for the Trust Game and for the Control game - performance for the tough Principal can be lower than for the generous one, despite the low performance is feasible for the performance set, proposed by the highly altruistic Principal.

The motivation crowding out is presented in the Core model as well, though it doesn't result in the crowding-out on performance, because the signaling of generosity goes through offering the performance level itself, not the (non-trivial) set of the performance levels, as in the Trust and Control games.

---

[32]The set of performance levels can be considered as a subset of the disobedience options: $\widehat{Q} \subset \overset{\circ}{Q}$. Alternatively, the set of disobedience options can be considered as a set of all possible Agent's actions.

In the Core model, the Agent attached to the generous Principal gets higher utility from an interaction. Indeed, he is intrinsically motivated for high performance (the participation threshold is high) but is asked for lower performance, and, consequently, gets positive utility. The Agent attached to the tough Principal gets positive utility only for the separating equilibria with $\beta \in (\beta_1, 1]$ - see Figure 6. For $\beta \in (\beta_2, \beta_1]$ the Agent performs at his participation threshold and gets zero utility, for $\beta \leq \beta_2$ the Agent gets negative utility in the pooling equilibria.

Another common feature of the Trust and the Control games is the effect of the strength of an extrinsic motivator. If it is weak, not using it can not signal the Principal's generosity. In fact, in the Control Game if the threshold $\overline{q}$ is small, it is not used by both types[33] and has no effect at the performance level. In the Trust Game, if the fine $\overline{f}$ is small, it is used by both types (the two lines $q^A$ and $\widehat{q}^A$ are close to each other so that pooling with both types imposing the fine emerges).[34]

If the available extrinsic motivator is strong enough but not too strong and Agent's reciprocity intensity is high enough, the generous Principal can signal her generosity by not using the extrinsic motivator. The tough Principal prefers to reveal her type and to crowd out the Agent's intrinsic motivation in this way because she gets better performance by using the extrinsic motivator and compensating in this way the crowding out of the intrinsic motivation.

Finally, if the available extrinsic motivator is very strong, both types can benefit form using it. The generous Principal crowds out the intrinsic motivation to the pooling level but gets compensated through high additional extrinsic motivation of the Agent.

# 5   Conclusion

My paper contributes to the literature on the theory of incentives which goes beyond the neoclassical framework and takes into account social components of preferences - reciprocity, fairness, social norms etc. In particular, I build a simple and intuitive model of reciprocal altruism and show its relevance for the behavioral patterns systematically observed both in the lab and in the field. The modeling approach can be used for the analysis of the workplace relation.

The model predicts that crowding-out in performance can emerge as the equilibrium outcome for some values of parameters of the model. The crowding-out in performance is the situation in which imposing an extrinsic incentive decreases the intrinsic motivation and leads to a lower performance - the phenomena well known in human resources management. It can also be the case (for some other values of parameters) that crowding-out in incentives doesn't result in crowding-out in performance.

This paper is among a few others studies aimed at enriching the theory of incentives by taking into account intrinsic motivation. Further research can be devoted to considering other information structures - Agents can differ in productivity, the organization can have more complicated structure than only

---

[33]Alternatively, it can be used by both types - see case I of proposition 5.

[34]Only the Control game experiment provides an evidence supporting the reciprocal altruism model. There were no variation of the size of the extrinsic incentive, i.e. the fine, in the Trust Game experiment.

one Principal and only one Agent. Considering other social components of preferences, supported by the evidence from lab and field (negative reciprocity, concerns for equity etc.) can also be relevant in building the theory of incentives.

Finally, I don't claim that monetary incentives are not important. On the contrary, it is well known that the incentive payments play an important role in creating incentives - see e.g. Bolton and Dewatripont [2005] or Prendergast [1999]. However, there is growing evidence that workers are motivated to exert effort not only by the incentive payment or other extrinsic motivators but also by the intrinsic motivation. On top of this, the interaction between intrinsic and extrinsic motivation can play an important role. The result of such interaction can be motivation crowding-out (or -in). Therefore, taking into account the intrinsic motivation in the labor contract models should give better understanding of the workplace relation. Intrinsic motivation and extrinsic incentives should be considered as complements rather than substitutes in the modeling.

# Appendix

## Proof of Lemma 5

According to (10), for any Best Response acceptance rule $a(\cdot)$ holds

$$a(q) \in \arg\max_{a \in [0,1]} \left(B - C(q) + \widehat{\alpha}(q)\beta q\right) a$$

The solution to this program is easy to find. If, for some $q$, $B - C(q) + \widehat{\alpha}(q)\beta q > 0$ then $a(q) = 1$, if $B - C(q) + \widehat{\alpha}(q)\beta q < 0$, then $a(q) = 0$ for the corresponding values of $q$, finally, if $B - C(q) + \widehat{\alpha}(q)\beta q = 0$, then any $a$ is a solution .

According to lemma 4, the participation constraint $B - C(q) + \widehat{\alpha}(q)\beta q > 0$ is equivalent to the threshold $q < q^0(\widehat{\alpha}(q), \beta)$. This gives the characterization of the Best Response rule.

Since $\mu(q) \in [0, 1]$, then $\alpha_L \leq \widehat{\alpha}(q) \leq \alpha_H$, consequently, according to monotonicity of the function $q^0(\alpha, \beta)$ with respect to its first argument (see lemma 4) we obtain

$$q^0(\alpha_L, \beta) \leq q^0(\widehat{\alpha}(q), \beta) \leq q^0(\alpha_H, \beta)$$

which can be rewritten as $q_L^0 \leq q^0(\widehat{\alpha}(q), \beta) \leq q_H^0$. This proofs (12).

## Proof of Lemma 6

*Proof.* Since $q_j^*$ $(j = L, H)$ are the elements of the Best Response of the Principal, the two inequalities hold:

$$(q_L^* - \alpha_L C(q_L^*))a^*(q_L^*) \geq (q_H^* - \alpha_L C(q_H^*))a^*(q_H^*)$$
$$(q_H^* - \alpha_H C(q_H^*))a^*(q_H^*) \geq (q_L^* - \alpha_H C(q_L^*))a^*(q_L^*)$$

Summing them up gives

$$(\alpha_H - \alpha_L)(C(q_L^*)a^*(q_L^*) - C(q_H^*)a^*(q_H^*)) \geq 0$$

which proves the first statement in the Lemma.

For the second claim of the Lemma, notice that if $a^*(q_H^*) = 1$, then

$$C(q_L^*) \geq C(q_L^*)a^*(q_L^*) \geq C(q_H^*)$$

which gives $q_L^* \geq q_H^*$. □

## Proof of Lemma 7

*Proof.* Let $q', q'' \in \text{supp } \sigma_H^*$ and $q', q'' \in \text{supp } \sigma_L^*$ are the two offers made by both types in an equilibrium. Then, the two types should be indifferent between the two offers:

$$(q' - \alpha_H C(q'))a^*(q') = (q'' - \alpha_H C(q''))a^*(q'')$$
$$(q' - \alpha_L C(q'))a^*(q') = (q'' - \alpha_L C(q''))a^*(q'')$$

This gives

$$\alpha_H(C(q'')a^*(q'') - C(q')a^*(q')) = \alpha_L(C(q'')a^*(q'') - C(q')a^*(q')) =$$
$$= q''a^*(q'') - q'a^*(q') \quad (28)$$

The first equality gives $C(q'')a^*(q'') - C(q')a^*(q') = 0$. This, in turn, gives

$$\frac{a^*(q'')}{a^*(q')} = \frac{q'}{q''} = \frac{C(q')}{C(q'')}$$

But the second part of this equality can't hold for a convex function $C(q)$ if $q' \neq q''$. This finishes the proof. □

## Proof of Lemma 8

*Proof.* According to Lemma 7 the pooling offer $q_p^*$ should be unique.

The belief should be consistent on the equilibrium path which means that $\mu^*(q_p^*) = \Pi$. This gives $\widehat{\alpha}(q_p^*) = \Pi$, according to (7). Thus, according to Lemma 5, the equilibrium acceptance rule $a^*(\cdot)$ should have a threshold $\widehat{q}(q_p^*) = q_E^0$. So, to have offer $q_p^*$ accepted, it is necessary to have $q_p^* \leq q_E^0$.

To prove sufficiency, note that any profile $(q_p^*, a^*(\cdot); \mu^*(\cdot))$ with offer $q_p^* \leq q_E^0$, acceptance rule $a^*(q_p^*) = 1$, $a^*(q) = 0$ for any other $q \neq q_p^*$ supported by beliefs $\mu^*(q_p^*) = \Pi$, $\mu^*(q) = 0$ for any other $q \neq q_p^*$ constitutes a PBE of the signaling game. □

## Proof of Lemma 9

*Proof.* I start the proof with two general properties of the equilibrium offers' supports.

**The 2-point support property.** The supports of equilibrium offer supp $\sigma_j^*$ contains at most two offers for each $j = L, H$.

To prove this, notice that any $q \in$ supp $\sigma_j^*$ should be Best Responses to the acceptance rule and solve the program

$$\max_{q \in Q} V(q, \alpha_j) a^*(q) \tag{29}$$

If $a^*(q) = 1$ for all $q \in$ supp $\sigma_j^*$ for some $j$, then all solutions of (29) should solve the program

$$\max_{q \in Q} V(q, \alpha_j) \tag{30}$$

However, it's clear that a solution of the latter program can consist of at most two points because of inverted-U shape of $V(q)$.

In case of $a^*(q) < 1$ for some equilibrium offer $q$, it is necessary that $q = q_L^0$, as follows from (33) and Lemma 5. Consider $j$ such that $q_L^0 \in$ supp $\sigma_j^*$. For all the others equilibrium offers $\tilde{q} \in$ supp $\sigma_j$ holds $\tilde{q} < q_L^0$ - see (31) and $a^*(\tilde{q}) = 1$ and, consequently $\tilde{q}$ solves (29) and (30) as in the previous case and there can be only two values $q = q'^{(j)}, q''^{(j)}$ which solve (30) because of inverted-U shape of $V(q)$.

It is important now to notice that $q'^{(j)} < q_L^0$, $q''^{(j)} > q_L^0$. This follows from (quasi-)convexity of $V(q, \cdot)$. In fact, we have

$$V(q'^{(j)}, \alpha_j) = V(q''^{(j)}, \alpha_j) = V(q, \alpha_j) a^*(q)$$

with $a^*(q) < 1$, which means that simultaneously hold $V(q, \alpha_j) > V(q'^{(j)}, \alpha_j)$ and $V(q, \alpha_j) > V(q''^{(j)}, \alpha_j)$. This, in turn, leads to $q'^{(j)} < q < q''^{(j)}$ because of convexity of $V(q)$

But $q''^{(j)}$ can't be an equilibrium offer since $q''^{(j)} > q_L^0$, consequently $a^*(q''^{(j)}) = 0$ according to Lemma 5. So, if $a^*\left(q_L^0\right) < 1$ and $q_L^0$ is an equilibrium offer for type $j$ then there can be at most one other equilibrium offer $q'^{(j)}$ for this $j$. This finishes the proof of the 2-points support property.

**The support bounds for the separating equilibria.**[35] For the separating equilibria holds

$$\text{supp } \sigma_H^* \subset [0, q_L^0], \quad \text{supp } \sigma_L^* \subset [0, q_L^0] \tag{31}$$

To prove this, notice that for any $q_L^*$ there should be $\mu^*(q_L^*) = 0$, consequently, to be accepted on the equilibrium path (at least with some probability, otherwise the $\alpha_L$-Principal's payoff will be zero), $q_L^*$ should satisfy $q_L^* \leq q_L^0$, according to Lemma 5. This inequality together with (13) gives[36]

$$q_H^* < q_L^* \leq q_L^0 < q_H^0 \tag{32}$$

This gives (31).

This means, in particular, that[37]

$$a^*(q_H^*) = 1 \quad \text{for any } q_H^* \tag{33}$$

since $q_H^* < q_H^0$, according to (32).

**The separating equilibria.**

The functions $V(q, \alpha_j)$ are increasing in $q$ for $q \in [0, q_j^P]$. For the case $q_L^0 \leq q_H^P$ this means that both functions $V(q, \alpha_j)$ are increasing in $q$ on supp $\sigma_j^*$, according to (31). So, if there are more than 1 offer $q \in \text{supp } \sigma_H^* \cup \text{supp } \sigma_L^*$ such that $a^*(q) = 1$ then both types must choose the highest of these offers to maximize payoff. Together with (33) this means that there is only one offer $\{q = q_H^*\} = \text{supp } \sigma_H^*$. Then, there should necessary be $q_L^* = q_L^0$ because only in this case it is possible to have $a^*(q_L^*) < 1$ - see Lemma 5 which leads to $\{q_L^0\} = \text{supp } \sigma_L^*$. This proofs the first part of the Lemma.

For the case of $q_L^0 > q_H^P$ the functions $V(q, \alpha_j)$ are no more monotone on $[0, q_L^0]$ and previous argument can't be used. So, the one-point support property doesn't hold but the 2-points support is guaranteed.

**The semi-separating equilibria.**

The 2-points support property together with the monotonicity Lemma 6 suggest that only the following three structures of the separating and pooling parts of the equilibria are possible:

$$q_H^* < q_p^* < q_L^*, \quad q_H^* < q_p^*, \quad q_p^* < q_L^*$$

However, the first one isn't possible for $q_L^0 \leq q_H^P$. To justify this, the same argument as for the separating equilibria can be applied: since the function $V(q, \alpha_H)$ is increasing on $[0, q_L^0]$ and $a^*(q_H^*) = a^*(q_p^*) = 1$ (it is possible to have $a^*(q_L^*) < 1$ only for the right-most point of the equilibrium offers), there is only one offer in the support of the $\alpha_H$-type: $\{q = q_H^*\} = \text{supp } \sigma_H^*$. For the other

---

[35]it can be shown the property holds for the separating parts of the semi-separating equilibria.

[36]It is possible that some elements of the equilibrium are not presented, then we can have, e.g. $q_H^* < q_p^* \leq q_L^0 < q_H^0$.

[37]Notice that we can guarantee that $a^*(q) = 1$ only for $q \in \text{supp } \sigma_H^* \setminus \text{supp } \sigma_L^*$.

two configurations $a^*(q) < 1$, where $q$ is the highest equilibrium offer, must hold to avoid deviation of the $\alpha_H$-type to this highest offer.

This finishes the proof of the necessity part, i.e. I have shown that other equilibrium structures are not possible. Existence is obvious. $\qquad\square$

## Proof of Proposition 2

*Proof.* **The case** $\beta > \beta_2 \Leftrightarrow q_H^P < q_L^0$.

The idea is that in this case the $\alpha_H$-type preferred offer $q_H^P$ is feasible, i.e. should be accepted for any reasonable acceptance rule. Even offers higher than $q_H^P$ are feasible (in particular, can be feasible $q_L^P$), and $\alpha_L$-type will be better off by proposing one of them. So, reasonably, $\alpha_H$-type should make only here preferred offer; $\alpha_L$-type should make here preferred offer $q_L^P$ if it is feasible and the highest feasible offer it $q_L^P$ is not feasible.

Put formally, consider an equilibrium, pooling or (semi-)separating, with

$$q_H^P \notin \operatorname{supp} \sigma_H^* \tag{34}$$

and apply the intuitive criterion (14) with $q = q_H^P, \alpha = \alpha_H$.

Notice that $q_H^P$ is always accepted according to all Best Response acceptance rules - see Lemma 5

$$a\left(q_H^P\right) = 1 \text{ for all } a \in BR\left(\mathcal{A}', q\right) \text{ for any } \mathcal{A}' \subseteq \mathcal{A}$$

This leads to

$$V\left(q_H^P, a, \alpha_H\right) = V\left(q_H^P, \alpha_H\right) a\left(q_H^P\right) = V\left(q_H^P, \alpha_H\right) \cdot 1 \text{ for all } a \in BR(\mathcal{A}', q)$$

Then $\min\limits_{a \in BR(\mathcal{A}', q)} V\left(q_H^P, a, \alpha_H\right)$ doesn't depend on $a$ and

$$\min_{a \in BR(\mathcal{A}', q)} V\left(q_H^P, a, \alpha_H\right) = V\left(q_H^P, \alpha_H\right)$$

which is the global maximum of $V(q, \alpha_H)$.

Consequently, $V_H^* < V\left(q_H^P, \alpha_H\right)$ and the intuitive criterion (14) can't hold given (34).

So, it is necessary for the refined equilibrium to have $q_H^P \in \operatorname{supp} \sigma_H^*$. Moreover, there can't be any other offer in $\operatorname{supp} \sigma_H^*$.

Consider now an equilibrium with

$$\min\left\{q_L^P, q_L^0\right\} \notin \operatorname{supp} \sigma_L^* \tag{35}$$

If $q_L^P < q_L^0$, apply the intuitive criterion (14) with $q = q_L^P, \alpha = \alpha_L$. The same argument as for $q_H^P, \alpha_H$ applies. This leads to $\operatorname{supp} \sigma_L^* = \{q_L^*\}$.

If $q_L^0 \leq q_L^P$ and $q_L^0 \notin supp\ \sigma_L^*$, apply (14) with $q = q_L^* + \varepsilon$, $\alpha = \alpha_L$. Again, the same argument applies[38] and leads to the conclusion that necessary $\operatorname{supp} \sigma_L^* = \{q_L^0\}$.

Finally, if $a^*\left(q_L^0\right) < 1$ for $q_L^0 \in \operatorname{supp} \sigma_L^*$, apply (14) with $q = q_L^0 - \varepsilon$, $\alpha = \alpha_L$. This offer is always reasonably accepted and there exists $\varepsilon$ small enough such that $V_L^* < V(q, a, \alpha_L)$ because $V_L^* = V(q_L^0, \alpha_L) a^*(q_L^0)$ and $V(q, a, \alpha_L) = V(q_L^0 - \varepsilon, \alpha_L) \cdot 1$. So, an equilibrium with $a^*\left(q_L^0\right) < 1$ can't pass the intuitive criterion.

---

[38]We can't guarantee that $a(q_L^0) = 1$ yet, so deviation to $q_L^0$ isn't necessarily profitable

**The case** $\beta \leq \beta_2 \ \Leftrightarrow \ q_L^0 \leq q_H^P$.

Notice that then $q_L^0 \leq q_L^P$ and both functions $V(q, \alpha_H)$ and $V(q, \alpha_L)$ are increasing on $[0, q_L^0]$ in this case.

If equilibrium is pooling with $q_p^* < q_L^0$, consider deviation to $q = q_p^* + \varepsilon$ with small $\varepsilon$ such that $q < q_L^0$. This offer $q$ is accepted for any reasonable acceptance rule and increases $V$ for both types, consequently (14) doesn't hold for such $q$ and any $\alpha$.

Consider pooling equilibria with $q_p^* \geq q_L^0$. According to Lemma 8, $q_p^* \leq q_E^0$ holds for any pooling PBE.

Let $q$ be a deviation such that $q < q_E^0$. Then

$$\max_{a \in BR(\mathcal{A}, q)} V(q, a, \alpha_j) = V(q, \alpha_j) \cdot 1$$

because at least for belief $\mu(q) = \Pi$ the offer $q$ should be accepted for any reasonable acceptance rule since $q < q_E^0$ - see Lemma 5.

Notice also that

$$\min_{a \in BR(\mathcal{A}', q)} V(q, a, \alpha_j) = 0 \text{ for } \mathcal{A}' = \{\alpha_L\} \text{ or } \mathcal{A}' = \mathcal{A}$$

$$\min_{a \in BR(\mathcal{A}', q)} V(q, a, \alpha_j) = V(q, \alpha_j) \text{ for } \mathcal{A}' = \{\alpha_H\}$$

because if belief are concentrated on the subset which contains $\alpha_L$, then belief $\mu(q) = 0$ is possible and any offer $q > q_L^0$ (which is the case here) is rejected for any reasonable acceptance rule, so the minimal value of $V$ is zero. On the other hand, if belief is concentrated on $\alpha_H$, i.e. $\mu(q) = 1$, the offer $q < q_H^0$ (which is the case since $q < q_E^0 < q_H^0$) is accepted.

So, the only possibility to have the intuitive criterion (14) violated is to have $\mathcal{A} \setminus J(q) = \{\alpha_H\}$ and $V_H^* < V(q, \alpha_H)$. In other words, $\alpha_H$-type should be "reasonably" revealed by deviation to $q$ (see the definition of the set $J(q)$) and this deviation should be profitable, i.e. $q$ should be closer to $q_H^P$, compared to the distance between $q_p^*$ and $q_H^P$.

Three cases are possible as illustrated by figure 14: 1) $q_p^* \leq q_H^P$; 2) $q_H^P < q_p^* \leq q_L^P$; 3) $q_p^* > q_L^P$.

For the case 1 the deviations to $q < q_p^*$ are not profitable; the deviations to $q > q_p^*$ are profitable (at least for $q < q_L^P$) but they are profitable for both types or only for $\alpha_L$-type. The deviations to $q > q_L^P$ can be profitable for $\alpha_L$-type only. Consequently, any deviation to $q > q_p^*$ can't reasonably reveal $\alpha_H$-type. So, all the equilibrium from area 1 pass the intuitive criterion.

For the case 2 the deviation to $q = q_H^P$ is profitable for $\alpha_H$-type and isn't profitable for the $\alpha_L$-type, so $\alpha_H$-type is revealed by such deviation and the equilibria from area 2 don't satisfy the intuitive criterion.

Finally, for the case 3, the revealing profitable deviation for the $\alpha_H$-type is constructed in the following way. Since $q_p^* > q_L^P$ there exists $q_L' < q_L^P$ (and then $q_L' < q_p^*$) such that $V(q_L', \alpha_L) = V(q_p^*, \alpha_L)$. Notice that $q_L'$ is on the increasing part of $V(q, \alpha_L)$, so any $q < q_L'$ is not a profitable deviation for $\alpha_L$. However, the deviation to $q_L'$ is profitable for $\alpha_H$-type (the proof is below), and then the deviation to $q = q_L' - \varepsilon$ is the revealing profitable deviation for $\alpha_H$ type.

Now I proof that the deviation to $q_L'$ is profitable for $\alpha_H$-type. Notice that the offers $q = q_L', q_p^*$ are accepted for any Best Response acceptance rule. So,

$$V(q, a, \alpha_j) = V(q, \alpha_j) \cdot 1 = q - \alpha_j C(q) + \alpha_j B$$
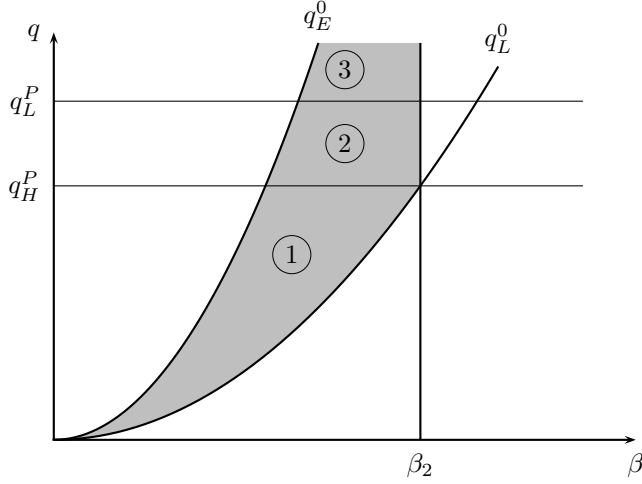
41

Figure 14:

Since $V(q'_L, \alpha_L) = V(q^*_p, \alpha_L)$, we have

$$q^*_p - \alpha_L C(q^*_p) = q'_L - \alpha_L C(q'_L)$$

By using this, we get

$$q'_L - \alpha_H C(q'_L) = q'_L - \alpha_L C(q'_L) + (\alpha_H - \alpha_L) C(q'_L) =$$
$$= q^*_p - \alpha_L C(q^*_p) + (\alpha_H - \alpha_L) C(q'_L) =$$
$$= q^*_p - \alpha_H C(q^*_p) + (\alpha_H - \alpha_L) C(q^*_p) + (\alpha_H - \alpha_L) C(q'_L) =$$
$$= q^*_p - \alpha_H C(q^*_p) + (\alpha_H - \alpha_L)(C(q^*_p) - C(q'_L))$$

The second term is positive, so

$$q'_L - \alpha_H C(q'_L) > q^*_p - \alpha_H C(q^*_p)$$

which gives the required inequality:

$$V(q'_L, \alpha_H) > V(q^*_p, \alpha_H)$$

This finishes the analysis of case 3 and the pooling equilibrium case. We've got that only the equilibria from area 1 satisfy the intuitive criterion.

If equilibrium is separating then, according to Lemma 9, $q^*_H < q^*_L = q^0_L$ and deviation to $q = q^*_H + \varepsilon$ doesn't satisfy (14) for $\alpha_H$ due to the same argument as for the pooling equilibrium with $q^*_p < q^0_L$.

If equilibrium is semi-separating, then for the structure $q^*_p < q^*_L = q^0_L$ the deviation to $q = q^*_p + \varepsilon$ for $\alpha_H$ doesn't satisfy (14) due to the same argument.

For the structure $q^*_H < q^*_p$, consider deviation of the separating part of equilibrium to $q$ which is closer to $q^P_H$: $q = q^*_H + \varepsilon$ if $q^*_H < q^P_H$, $q = q^*_H - \varepsilon$ if $q^*_H > q^P_H$ (it's impossible to have $q^*_H = q^P_H$ because then $\alpha_H$-type strongly prefers $q^P_H$ to the pooling part of the equilibrium candidate $q^*_p$, which can't be

the case in equilibrium). First, $J(q) = \{\alpha_L\}$ because the $\alpha_L$ type is concentrated on the pooling part of the equilibrium $q_p^*$; since $q < q_p^*$ and $V(q, \alpha_L)$ is increasing, deviation to $q$ is unprofitable for $\alpha_L$-type. Second, in the intuitive criterion (14) $a \in BR(\{\alpha_H\}, q)$ which means that offer $q$ is accepted with probability 1: $a(q) = 1$ and then $V_H^* = V(q_H^*, \alpha_H) < \min\limits_{a \in BR(\{\alpha_H\}, q)} V(q, a, \alpha_H)$ with $V(q, a, \alpha_H) \equiv V(q, \alpha_H)$).

$\square$

## Proof of Proposition 3

*Proof.* For the case 1 all the Principals have $\alpha > \alpha^\times$, then, according to Claim 1, the preferred output for all types is feasible. So, accepting all offers $q^P(\alpha)$ is reasonable Best Response for the Agent and gives to all types of the Principal their unconstrained maximal utility.

For the case 2, first proof that there exists unique solution to (16).

For $\alpha = \alpha^\times$ we have $E_{\alpha^\times}[\alpha] < \alpha^\times$, and, consequently $q^0(E_{\alpha^\times}[\alpha]) < q^0(\alpha^\times) = q^P(\alpha^\times) < q^P(E_{\alpha^\times}[\alpha])$, so the left-hand side of (16) is smaller than than the right-hand side.

For $\alpha = \alpha_2$ holds $E_{\alpha^\times}[\alpha] = E\alpha$. So, (15) means that the left-hand side of (16) is more than the right-hand side.

Since both sides of (16) are continuous, there exists unique solution to this equation.

Second, consider the Agent's acceptance rule given the Principal's offer. Clearly, all offers $q > \tilde{q}^0$ should be reasonable rejected, offers $q < \tilde{q}^0$ should be reasonably accepted with probability 1. The offer $q = \tilde{q}^0$ should be accepted with probability 1 on the equilibrium path because accepting it with probability $a < 1$ will make Principal's deviation to $\tilde{q}^0 - \varepsilon$ profitable.

Finally, ant type of the Principal can't do better since those with $\alpha > \tilde{\alpha}$ implement their preferred output, and those with $\alpha < \tilde{\alpha}$ could do better only by implementing $q > \tilde{q}^0$ which are rejected.

Case 3 is considered in the same way as the pooling part ($\alpha < \tilde{\alpha}$) in case 2.

$\square$

## Proof of Claim 2

*Proof.* The statements in 1 and 2 are trivial since the Agent has full flexibility and hence chooses his preferred back-transfer.

For the third point, notice that if the desired back-transfer $\hat{q} \le q^A$, then paying back $q^A$ will not impose fine and will maximize the Agent's utility ($q^A$ is the global maximizer).

Consider the case $\hat{q} > q^A$.

Notice that $\tilde{q}^A(\hat{\alpha})$ is constructed in such way that

$$\overset{\circ}{U}(q) > \overset{\circ}{U}(q^A) - f \quad \text{for} \quad q^A(\hat{\alpha}, \beta) < \hat{q} < \tilde{q}^A(\hat{\alpha}, \beta) \tag{36}$$

$$\overset{\circ}{U}(q) < \overset{\circ}{U}(q^A) - f \quad \text{for} \quad \hat{q} > \tilde{q}^A(\hat{\alpha}, \beta) \tag{37}$$

where $\overset{\circ}{U}(q)$ is the Agent's utility without taking into account the possibility of fine: $U(q) = \overset{\circ}{U}(q) - f I_{q < \hat{q}}$.

In (36) the Agent prefers to diverge from $q^A$ to $q > q^A$ as such divergence isn't too high whereas in (37) the Agents prefers to pay fine. $\square$

## Proof of the Proposition 4

*Proof.* The optimality of the Principal's decision given the Agents's belief is evident from Claim 2. We should check the incentives compatibility conditions and the crowding-out condition.

Consider the case $q_{LH} \leq \widetilde{q}_{LL}$.

The Principal's incentive compatibility constraints are

$$\pi_H(q_{HH} - \alpha_H C(q_{HH})) + (1 - \pi_H) \cdot 0 \geq \widetilde{q}_{LL} - \alpha_H C(\widetilde{q}_{LL})$$
$$\widetilde{q}_{LL} - \alpha_L C(\widetilde{q}_{LL}) \geq \pi_L(q_{HH} - \alpha_L C(q_{HH})) + (1 - \pi_L) \cdot 0$$

which are equivalent to (22) and (23), correspondingly.

The inequality $\widehat{\pi}_H \leq 1$ holds since the denominator in (22) is positive and the inequality is then equivalent to

$$\frac{C(q_{HH}) - C(\widetilde{q}_{LL})}{q_{HH} - \widetilde{q}_{LL}} \leq \frac{1}{\alpha_H}$$

The left-hand side $\frac{1}{\alpha_H} \geq 1$. The right-hand side is the slope of the secant line to the graph of the convex function $C(q)$ between the points with $q = \widetilde{q}_{LL}$ and $q = q_{HH}$, which is smaller than the slope of the tangent line at the point with $q = q_{HH}$ which is equal to $C'(q_{HH}) = \alpha_H \beta_H < 1$. So, the inequality holds.

The inequality $\widehat{\pi}_L > 0$ holds since both the numerator and the denominator of the fraction are positive.

Finally, we need to check the crowding-out condition $\pi q_{HH} \geq \widetilde{q}_{LL}$, where $\pi$ is the objective probability of the selfish Agents. Since the selfish Principal has belief $\pi_L$ which is biased downward, it is sufficient to prove that $\widehat{\pi}_L q_{HH} \geq \widetilde{q}_{LL}$. Substituting $\widehat{\pi}_L$ into the inequality, we get

$$\frac{\widetilde{q}_{LL} - \alpha_L C(\widetilde{q}_{LL})}{q_{HH} - \alpha_L C(q_{HH})} q_{HH} \geq \widetilde{q}_{LL}$$

Since the denominator is positive, this inequality is equivalent to

$$\alpha_L \widetilde{q}_{LL} q_{HH} \left( \frac{C(q_{HH})}{q_{HH}} - \frac{C(\widetilde{q}_{LL})}{\widetilde{q}_{LL}} \right) \geq 0$$

which holds since $q_{HH} > \widetilde{q}_{LL}$, $\widetilde{q}_{LL} < \widetilde{q}_{LH}$ and it's assumed that $\widetilde{q}_{LH} < q_{HH}$. $\square$

## Proof of the Corollary 2

*Proof.* The condition $q_{LH} \leq \widetilde{q}_{LL}$ is equivalent to $C(q_{LH}) \leq C(\widetilde{q}_{LL})$. Since $C(\widetilde{q}_{LL}) = f$, it leads to $C(q_{LH}) \leq f$, so that $f_1 = C(q_{LH})$.

Now check the condition $\widetilde{q}_{LH} \leq q_{HH}$.

The back-transfer $\widetilde{q}_{LH}$ is determined, according to Claim 2 by

$$\alpha_L \beta_H q_{LH} - C(q_{LH}) - f = \alpha_L \beta_H \widetilde{q}_{LH} - C(\widetilde{q}_{LH})$$

where $\widetilde{q}_{LH}$ is chosen on the decreasing part of the graph of the function $F(q) = \alpha_L \beta_H q - C(q)$ (see Figure 10). Consequently, $\widetilde{q}_{LH} \leq q_{HH}$ is equivalent to

$$\alpha_L \beta_H q_{LH} - C(q_{LH}) - f \geq \alpha_L \beta_H q_{HH} - C(q_{HH})$$

which can be rewritten as

$$f_2 \equiv (\alpha_L \beta_H q_{LH} - C(q_{LH})) - (\alpha_L \beta_H q_{HH} - C(q_{HH})) = f_2 \geq f$$

Finally, to make sure that the interval $[f_1, f_2]$ is non-empty, we should check that $f_1 \leq f_2$. This leads to

$$\alpha_L \alpha_H \beta_H^2 \left( \frac{q_{LH}}{\alpha_L \beta_H} - \frac{q_{HH}}{\alpha_H \beta_H} \right) \leq C(q_{LH})$$

for a generic cost function.

For the quadratic cost function $C(q) = \frac{c}{2}q^2$, taking into account that $q_{ij}$ are determined by $C'(q_{ij}) = \alpha_i \beta_j$, and substituting this into the last inequality, one can check that the left-hand side is equal to zero, so that the inequality always holds.

Finally, for given $\alpha_L, \alpha_H, \beta_H$, and $f \in [f_1, f_2]$, one can obtain the threshold values $\widehat{\pi}_L \geq 0$, $\widehat{\pi}_H \leq 1$ from (22) and (23), and take the values $\pi_L$ and $\pi_H$, satisfying $\pi_H \geq \widehat{\pi}_H$, $\pi_L \leq \widehat{\pi}_L$. For these parameters, according to Proposition 4, the equilibrium of the signaling game is the separating crowding-out equilibrium. $\qquad\square$

## Proof of Proposition 5

*Proof.* For the case $\overline{q} \geq q_{LH}$, the incentives compatibility constraints for the Principal are:

$$\pi_H(q_{HH} - \alpha_H C(q_{HH})) \geq \overline{q} - \alpha_H C(\overline{q})$$
$$\overline{q} - \alpha_L C(\overline{q}) \geq \pi_L(q_{HH} - \alpha_L C(q_{HH}))$$

These constraints are equivalent to the conditions $\pi_L \leq \widehat{\pi}_L$ and $\pi_H \geq \widehat{\pi}_H$, which are assumed to hold.

Check the crowding-out condition $\pi q_{HH} \geq \overline{q}$. Since $\pi > \pi_L$, the inequality $\pi_L q_{HH} \geq \overline{q}$ is stronger. I check the latter inequality for $\pi_L = \widehat{\pi}_L$.

Substituting the formulae for $\widehat{\pi}_L$, we get

$$\frac{\widetilde{q} - \alpha_L C(\widetilde{q})}{q_{HH} - \alpha_L C(q_{HH})} q_{HH} \geq \widetilde{q}$$

after rearranging it leads to

$$\frac{C(\widetilde{q})}{\widetilde{q}} \leq \frac{C(q_{HH})}{q_{HH}}$$

which is equivalent to $\widetilde{q} \leq q_{HH}$ since the function $C(q)$ is convex. The latter inequality is assumed to hold.

So, at least for $\pi_L$ close to $\widehat{\pi}_L$ the crowding-out condition holds.

For the case of $\overline{q} < q_{LH}$, the incentives compatibility constraints for the Principal are:

$$\pi_H(q_{HH} - \alpha_H C(q_{HH})) \geq \pi_H(q_{LH} - \alpha_H C(q_{LH})) + (1 - \pi_H)(\overline{q} - \alpha_H C(\overline{q}))$$
$$\pi_L(q_{LH} - \alpha_L C(q_{LH})) + (1 - \pi_L)(\overline{q} - \alpha_L C(\overline{q})) \geq \pi_L(q_{HH} - \alpha_L C(q_{HH}))$$

As in the previous case, the crowding-out condition will hold at least for $\pi_L$ close to $\widetilde{q}_L$ if $\widehat{\pi}_L q_{HH} \geq \overline{q}$. Substituting the expression for $\widehat{\pi}_L$ gives

$$\frac{\overline{q} - \alpha_L C(\overline{q})}{[q_{HH} - \alpha_L C(q_{HH})] + [\overline{q} - \alpha_L C(\overline{q})] - [q_{LH} - \alpha_L C(q_{LH})]} q_{HH} \geq \overline{q}$$

which can be rearranged to

$$\alpha_L q_{HH} \left[ \frac{C(q_{HH})}{q_{HH}} - \frac{C(\overline{q})}{\overline{q}} \right] \geq (q_{LH} - \overline{q}) \left[ \alpha_L \frac{C(q_{LH}) - C(\overline{q})}{q_{LH} - \overline{q}} - 1 \right] \qquad (38)$$

The left-hand side term $\frac{C(q_{HH})}{q_{HH}} - \frac{C(\overline{q})}{\overline{q}} > 0$ because $q_{HH} > q_{LH} > \overline{q}$.

The right-hand side term $\frac{C(q_{LH}) - C(\overline{q})}{q_{LH} - \overline{q}} < 1$, because it's a slope of the secant line to the graph of the increasing convex function $C(q)$, which is lower than the slope of the tangent line at the right edge of the interval $[\overline{q}, q_{LH}]$, $C'(q_{LH})$, for which we have $C'(q_{LH}) = \alpha_L \beta_H < 1$.

So, the right-hand side in (38) is positive, the left-hand side is negative, and, consequently, the inequality (38) holds.

$\square$

## Proof of Proposition 6

*Proof.* The proposition is established by checking the equilibrium conditions case by case.

Consider the no-control pooling equilibrium candidate. Each of the Principals shouldn't have an incentive to deviate to control, in which case the output $\overline{q}$ will be obtained. For the pooling equilibrium the Agent's belief on the Principal's type on the equilibrium path is $E\alpha$, and the pro-social Agent will perform at the level $q_{EH}$, determined by $C'(q_{EH}) = E\alpha\beta_H$, the selfish Agent will perform at level $q = 0$. So, the two Best Response conditions for the two types of Principal are

$$V_H = \pi_H(q_{EH} - \alpha_H C(q_{EH})) \geq \overline{q} - \alpha_H C(\overline{q})$$
$$V_L = \pi_L(q_{EH} - \alpha_L C(q_{EH})) \geq \overline{q} - \alpha_L C(\overline{q})$$

The two inequalities hold for small $\overline{q}$, since the right-hand sides are equal to 0 for $\overline{q} = 0$. The first condition which becomes binding for small $\overline{q}$ determines the threshold $q_1$.

For the large $\overline{q}$, the right-hand sides of the two inequalities are negative. By decreasing $\overline{q}$, the inequality for $V_L$ becomes bonding and determines the threshold $q_6$.

Consider the control pooling equilibrium. The Principals' Best Response conditions are (the Agent will reasonable believe that the Principal deviating to no-control should be the selfish one)

$$V_H = \overline{q} - \alpha_H C(\overline{q}) \geq \pi_H(q_{HH} - \alpha_H C(q_{HH}))$$
$$V_L = \overline{q} - \alpha_L C(\overline{q}) \geq \pi_L(q_{HH} - \alpha_L C(q_{HH}))$$

The inequality for the $V_H$ is stronger and determines the lower and upper bounds for the values of $\overline{q} \in [q_3, q_4]$ for which the control pooling equilibrium emerges.

The case of the separating equilibrium is partially considered in Proposition 5. The conditions for the separating equilibrium to the right of the control pooling region are the same, and the "right-hand side" separating equilibrium emerges due to non-monotonicity of the payoff functions. The regions for $\overline{q}$ are $[q_2, q_3]$ on the left and $[q_4, q_5]$ on the right.

So, all the possible pure strategies equilibria are considered. The regions of the values of $\overline{q}$ not covered by the pure strategies equilibria, should bring the mixed strategies equilibria.

$\square$

# References

Klaus Abbink, Bernd Irlenbusch, and Elke Renner. The moonlighting game An experimental study on reciprocity and retribution. *Journal of Economic Behavior and Organization*, 42(2):265–277, 2000.

James Andreoni. Philanthropy, Handbook of Giving, Reciprocity and Altruism, SC. Kolm and J. Mercier Ythier, eds, 2006.

Charles Bellemare and Bruce Shearer. Gift Exchange within a Firm: Evidence from a Field Experiment. *Cahier de recherche/Working Paper CIRPÉE*, 7:08, 2007.

Roland Bénabou and Jean Tirole. Incentives and prosocial behavior. *American Economic Review*, 96(5):1652–1678, 2006a.

Roland Bénabou and Jean Tirole. Belief in a Just World and Redistributive Politics. *Quarterly Journal of Economics*, 121(2):699–746, 2006b.

Roland Bénabou and Jean Tirole. Intrinsic and Extrinsic Motivation. *Review of Economic Studies*, 70(3):489–520, 2003.

Joyce Berg, John Dickhaut, and Kevin McCabe. Trust, Reciprocity, and Social History. *Games and Economic Behavior*, 10(1):122–142, 1995.

Sandra H. Berry and David E. Kanouse. Physician Response to a Mailed Survey an Experiment in Timing of Payment. *Public Opinion Quarterly*, 51(1):102–114, 1987.

Gary E. Bolton and Axel Ockenfels. Does Laboratory Trading Mirror Behavior in Real World Markets. *Fair Bargaining and Competitive Bidding on EBay, March*, 2008.

Patrick Bolton and Mathias Dewatripont. *Contract Theory*. MIT Press, 2005.

Gary Charness, Ernan Haruvy, and Doron Sonsino. Social distance and reciprocity: An Internet experiment. *Journal of Economic Behavior and Organization*, 63(1):88–103, 2007.

In-Koo Cho and David M. Kreps. Signaling Games and Stable Equilibria. *Quarterly Journal of Economics*, 102(2):179–221, 1987.

Martin Dufwenberg and Uri Gneezy. Measuring Beliefs in an Experimental Lost Wallet Game. *Games and Economic Behavior*, 30(2):163–182, 2000.

Robert Dur. Gift Exchange in the Workplace: Money or Attention? *Tinbergen Institute Discussion paper*, 2008.

Tore Ellingsen and Magnus Johannesson. Pride and prejudice: The human side of incentive theory. *American Economic Review*, 98(3):990–1008, June 2008.

Armin Falk. Gift Exchange in the Field. *Econometrica*, 75(5):1501–1511, 2007.

Armin Falk and Urs Fischbacher. A theory of reciprocity. *Games and Economic Behavior*, 54(2):293–315, 2006.

Armin Falk and Michael Kosfeld. The Hidden Costs of Control. *American Economic Review*, 96(5):1611–1630, 2006.

Ernst Fehr and Armin Falk. Wage Rigidity in a Competitive Incomplete Contract Market. *Journal of Political Economy*, 107(1):106–134, 1999.

Ernst Fehr and Urs Fischbacher. Third-party punishment and social norms. *Evolution and Human Behavior*, 25(2):63–87, 2004.

Ernst Fehr and John A. List. The Hidden Costs and Returns of Incentives-Trust and Trustworthiness Among CEOs. *Journal of the European Economic Association*, 2(5):743–771, 2004.

Ernst Fehr and Bettina Rockenbach. Detrimental effects of sanctions on human altruism. *Nature*, 422:137–140, 2003.

Ernst Fehr, George Kirchsteiger, and Arno Riedl. Does Fairness Prevent Market Clearing? An Experimental Investigation. *Quarterly Journal of Economics*, 108(2):437–460, 1993.

Ernst Fehr, Simon Gächter, and George Kirchsteiger. Reciprocity as a Contract Enforcement Device: Experimental Evidence. *Econometrica*, 65(4):833–860, 1997.

Drew Fudenberg and Jean Tirole. *Game Theory*. Mit Press, 1991.

Simon Gächter and Armin Falk. Reputation and Reciprocity: Consequences for the Labour Relation. *Scandinavian Journal of Economics*, 104(1):1–26, 2002.

Uri Gneezy. Does high wage lead to high profits? An experimental study of reciprocity using real effort. *Graduate School of Business, University of Chicago*, 2002.

Uri Gneezy and John A. List. Putting Behavioral Economics to Work: Testing for Gift Exchange in Labor Markets Using Field Experiments. *Econometrica*, 74(5):1365–1384, 2006.

Werner Güth, Rolf Schmittberger, and Bernd Schwarze. An experimental analysis of ultimatum bargaining. *Journal of Economic Behavior and Organization*, 3(4):367–88, 1982.

Heike Hennig-Schmidt, Bettina Rockenbach, and Abdolkarim Sadrieh. In search of workers real effort reciprocity–A field and a laboratory experiment. *Governance and the Efficiency of economic SYstems DP*, 55, 2005.

Daniel Kahneman, Jack L. Knetsch, and Richard Thaler. Fairness as a Constraint on Profit Seeking: Entitlements in the Market. *American Economic Review*, 76(4):728–741, 1986.

Sebastian Kube, Michel André Maréchal, and Clemens Puppe. Putting Reciprocity to Work - Positive Versus Negative Responses in the Field. *SSRN eLibrary*, 2006.

David K. Levine. Modeling Altruism and Spitefulness in Experiments. *Review of Economic Dynamics*, 1(3):593–622, 1998.

49

John A. List. Field Experiments: A Bridge Between Lab and Naturally-Occurring Data. *NBER Working Paper*, 2007.

John A. List and Steven D. Levitt. What do laboratory experiments tell us about the real world. *NBER Working Paper*, 2005.

George Loewenstein, Ted O'Donoghue, and Matthew Rabin. Projection Bias In Predicting Future Utility. *Quarterly Journal of Economics*, 118(4):1209–1248, 2003.

Michel A. Maréchal and Christian Thöni. Do Managers Reciprocate? Field Experimental Evidence from a Competitive Market. *SSRN eLibrary*, 2007.

Harry J. Paarsch and Bruce S. Shearer. The Response to Incentives and Contractual Efficiency: Evidence from a Field Experiment. *Cahier de recherche CIRPE/Working Paper*, 7:01, 2007.

Canice Prendergast. The Provision of Incentives in Firms. *Journal of Economic Literature*, 37(1):7–63, 1999.

Matthew Rabin. Incorporating Fairness into Game Theory and Economics. *American Economic Review*, 83(5):1281–1302, 1993.

Bruce Shearer. Piece Rates, Fixed Wages and Incentives: Evidence from a Field Experiment. *Review of Economic Studies*, 71(2):513–534, 2004.

Dirk Sliwka. Trust as a signal of a social norm and the hidden costs of incentive schemes. *American Economic Review*, 97(3):999–1012, 2007.

Jean Tirole. Rational irrationality: Some economics of self-management. *European Economic Review*, 46(4-5):633–655, 2002.